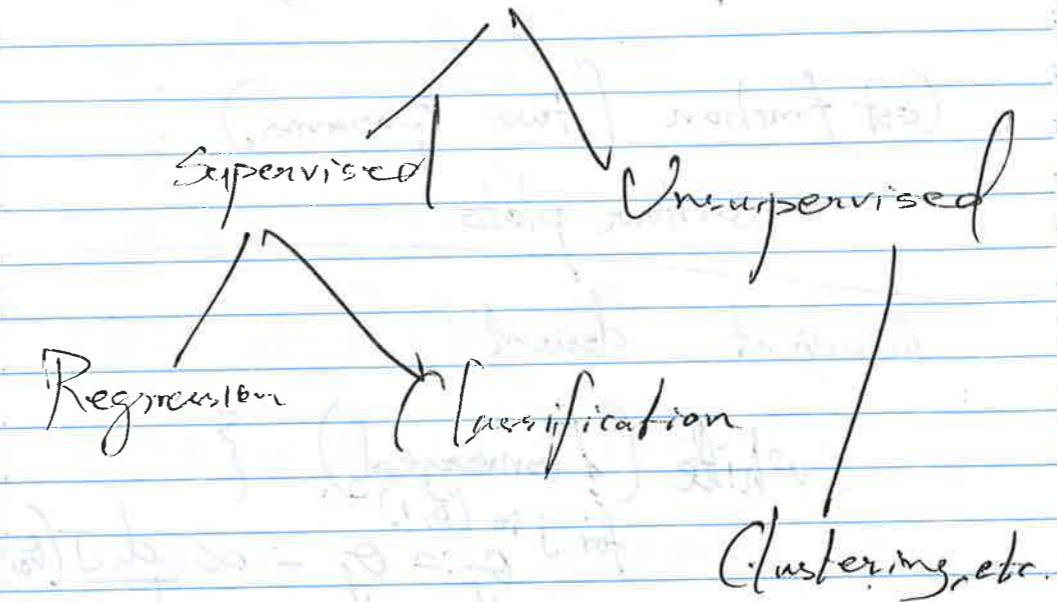


ML Coursera:

Week 1:

Linear regression



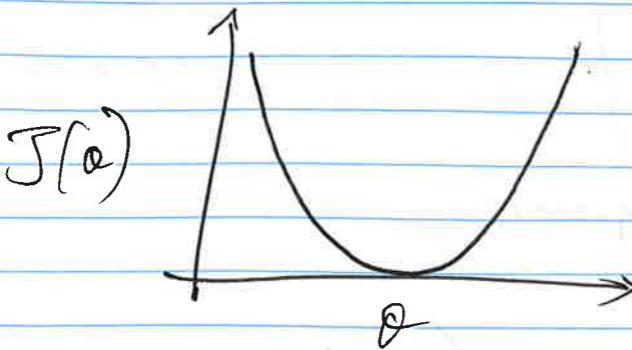
Linear regression (one variable):

$$h(x) = \theta_0 + \theta_1 x$$

Cost function:

$$\text{Minimize}_{\theta_0, \theta_1} J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h(x_i) - y_i)^2$$

Cost function (one param) :



Cost function (two params) :

- contour plots

Gradient descent :

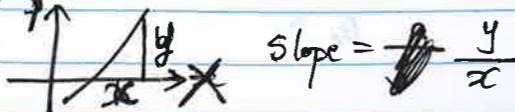
while (! converged) {
for j in (0, 1):
 $\theta_j := \theta_j - \alpha \frac{d}{d\theta_j} J(\theta)$ }

}
 α - learning rate

Simultaneous update

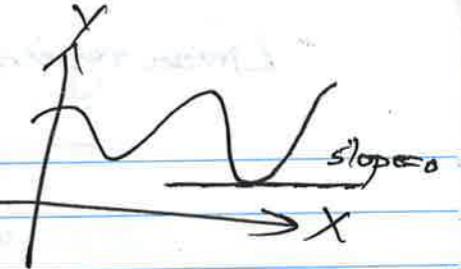
$\frac{d}{d\theta} (J(\theta)) \rightarrow$ Change in cost for
for a small change in θ

Slope of the tangent



At local minimum:

(one param)



$$\theta := \theta - \alpha \frac{d}{d\theta} (J(\theta))$$

say $J(\theta) = \text{local}$

At local minimum, $\theta = \theta_*$.

$$\begin{aligned} \theta &:= \theta - \alpha \frac{d}{d\theta} (\theta_* \alpha) \\ \theta &:= \theta - \alpha \alpha \\ \theta &:= \theta \end{aligned}$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (\theta x - y)^2 \quad (\text{here, } h(x) = \theta x)$$

$$\theta := \theta - \alpha \frac{d}{d\theta} \left(\frac{1}{2m} \sum_{i=1}^m (\theta x - y)^2 \right)$$

$$\theta := \theta - \alpha \frac{x^2}{2m} \sum_{i=1}^m \left(\frac{d}{d\theta} (\theta x - y) \right)$$

$$\theta := \theta + \alpha \frac{1}{m} \sum_{i=1}^m x (x - y)$$

$x = 0$ at local minimum

$$\theta := \theta$$



$$\theta := \theta - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (\theta x - y) \cdot x \quad (\text{Chain rule})$$

At local minima, θx is as close to y as possible.
Ideal case, $\theta x = y$.

Linear regression cost function:

- Convex function

↙ Only one ~~local~~.

No local minima, only global minimum

Batch gradient

Batch gradient descent

↳ Look at all training examples during each iteration

$$\left(\sum_{i=1}^m (h_\theta(x_i) - y_i) \right)$$

Other scenario - look at subset of training data in each iter

Gradient descent for linear regression:

$$h_\theta(x) = \theta_0 + \theta_1 x \quad (\text{Hypothesis / model})$$

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2 \quad (\text{Cost function})$$

$$\theta_0 := \theta_0 - \alpha \frac{d}{d\theta_0} \left[\frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2 \right]$$

$$\theta_0 := \theta_0 - \alpha \cdot \frac{1}{m} \cdot \sum_{i=1}^m (h_\theta(x_i) - y_i) \quad (\text{Chain rule})$$

$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} \left[\frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2 \right]$$

$$\theta_1 := \theta_1 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_\theta(x_i) - y_i) \cdot x_i \quad (\text{Chain rule})$$

WEEK 2:

Multivariate Linear Regression:

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 + \dots + \theta_n x_n$$

$$\rightarrow h_\theta(x) = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad (\text{Here, } x_0 = 1)$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\boxed{h_\theta(x) = \theta^T x}$$

$$\text{Cost fn.: } J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x_i) - y_i)^2$$

Gradient descent:

Repeat {

$$\theta_j := \theta_j - \alpha \frac{d}{d\theta_j} J(\theta)$$

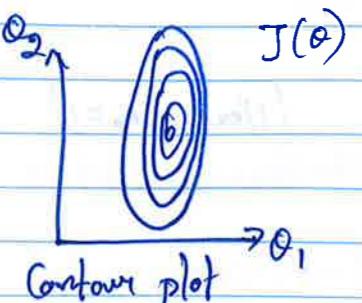
{Simultaneously update all
 θ_j for
 $j = 0, 1, \dots, n$ }

$$\theta_j := \theta_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_\theta(x_i) - y_i) \cdot x_i^j$$

Feature scaling:

→ Scale features to make gradient descent converge faster

(eg) $x_1 = \text{size}(0-2000 \text{ feet}^2)$
 $x_2 = \text{number of bedrooms } (1-5)$



This makes convergence slow

Scale features: $x_1 = \frac{\text{size}(\text{feet}^2)}{2000}$

$x_2 = \text{no. of bedrooms}$

$0 \leq x_1 \leq 1 ; 0 \leq x_2 \leq 1$



Mean normalization:

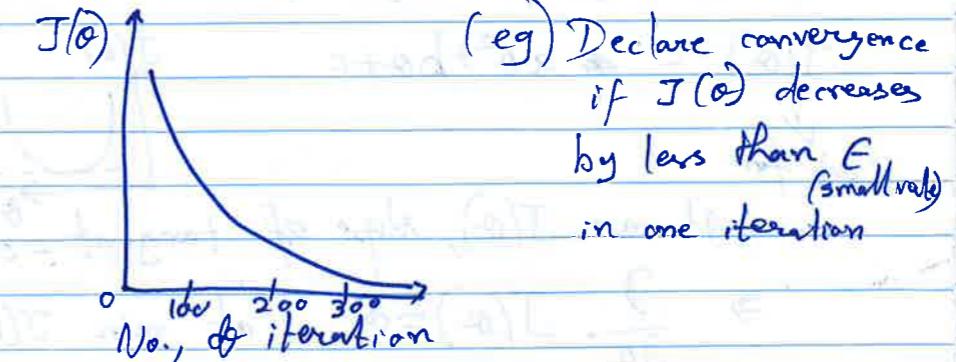
Replace x_i with $\frac{x_i - \mu_i}{s_i}$

(μ_i → mean of all the feature value over all training examples)

In general, $x_i \leftarrow \frac{x_i - \mu_i}{s_i}$ ($s_i \rightarrow \text{Range or std deviation}$)

Range will have the full impact of outliers, std deviation handles outliers better

Gradient descent - Learning rate:



(eg) Declare convergence if $J(\theta)$ decreases by less than E (E small rate) in one iteration

If $J(\theta)$ is increasing, choose smaller value for α (Learning rate)

→ α too small: slow convergence

→ α too large: may not converge.

Features and polynomial regression:

~~Quadratic model~~

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^3$$

(Cubic model)

↓ Make this cubic function a linear model by choosing features as

Features $\begin{cases} x_1 = x \\ x_2 = x^2 \\ x_3 = x^3 \end{cases}$ → Parts of cubic fn.

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

→ Now apply linear regression
 (Feature scaling is extremely important in this case)

Normal equation:

solve for θ analytically w/o iterations

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (\theta^T x_i - y_i)^2$$

(cost fn.)

At min $J(\theta)$, slope of tangent $= 0$

$$\Rightarrow \frac{\partial}{\partial \theta} J(\theta) = 0 \text{ at min } J(\theta)$$

Solve for θ .

This is for $\theta \rightarrow 0$.

For $\theta \in \mathbb{R}^{n+1}$

$$J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (\theta^T x_i - y_i)^2$$

$$\text{Set: } \frac{\partial}{\partial \theta_j} J(\theta) = 0 \text{ (for every } j)$$

Solve for $\theta_0, \theta_1, \dots, \theta_n$

Matrix of all training data

$$X = \begin{bmatrix} & & \\ & & \end{bmatrix} \quad \text{training examples}$$

$$y = \begin{bmatrix} & \\ & \end{bmatrix} \quad \text{m dimensional vector}$$

$$\theta = (X^T X)^{-1} X^T y$$

$$X \theta = Y$$

For X to have inverse, X has to be a square matrix.
So $\theta = X^{-1} Y$ may not always be possible.

$$\therefore X^T X \theta = X^T Y$$

$X^T X$ is a square matrix

$$\Rightarrow \theta = (X^T X)^{-1} X^T Y$$

Normal equation

Feature scaling isn't necessary for normal equation method

(since no iterations, or gradient descent, or learning rate involved).

With gradient descent, w/o feature scaling, a value of α (learning rate) that's applicable for one feature, may not apply for other features.

Gradient Descent

→ Need to choose α

→ Needs many iterations

→ Works well when no. of features (n) is large.

Normal equation

→ No α

→ No iterations

→ Expensive if n is large.

$(X^T X)^{-1}$ → $n \times n$ matrix

→ Large time complexity

What if $X^T X$ is non-invertible?

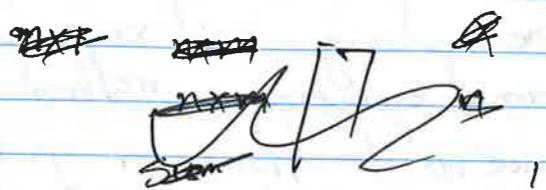
Could be due to:

↳ Redundant (linearly dependant) features

↳ Too many features ($m \leq n$)

Delete some features

Regularization



Vectorization:

Gradient descent:

$$\theta = \theta - \alpha \cdot \frac{1}{m} \cdot \sum ((\theta^T x - y) \cdot x)$$

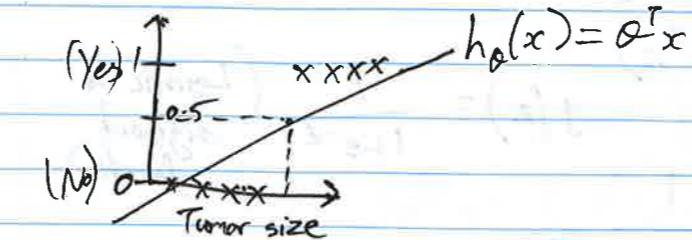
$$97 \quad \begin{pmatrix} 97 \\ + \end{pmatrix}$$

$$\theta = \theta - \alpha \cdot \frac{1}{m} \cdot \sum ((\theta^T x - y) \cdot x)^T$$

~~Def: Invertible~~
 ~~$X^T X$~~

Week 3:

Classification:

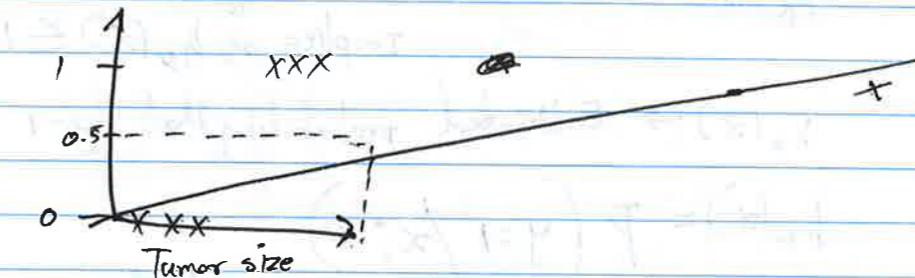


Apply Linear regression \rightarrow threshold classifier

If $h_\theta(x) \geq 0.5$, predict "y=1"

If $h_\theta(x) < 0.5$, predict "y=0"

But, outliers change the regression model:



One outlier example can have a big change in model.

Also, linear regression can result in

~~h_θ(x) > 1 or h_θ(x) < 0,~~

but in classification, y=0 or 1.

\Rightarrow Linear regression isn't applicable for classification problems.

Logistic regression: \rightarrow Want $0 \leq h_\theta(x) \leq 1$

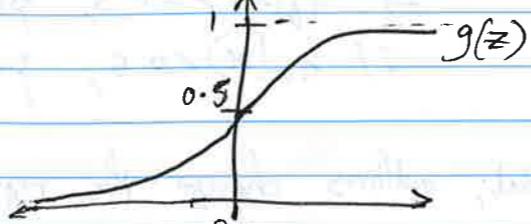
$$h_\theta(x) = g(\theta^T x)$$

$$h_\theta(x) = g(\theta^T x)$$

$$g(z) = \frac{1}{1+e^{-z}} \quad (\text{Logistic or sigmoid function})$$

$$h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$$

Sigmoid fn.:



$h_\theta(x)$

Implies $0 \leq h_\theta(x) \leq 1$

$h_\theta(x) \rightarrow$ Estimated probability that $y=1$ for x .

$$h_\theta(x) = P(y=1/x; \theta)$$

(probability that $y=1$ given x , parametrized by θ)

$$\text{P}(y=0/x; \theta) = 1 - P(y=1/x; \theta)$$

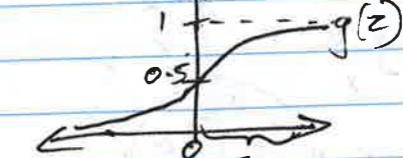
When $z=0$, $g(z) = 0.5$

~~$z=\infty, g(z)=1$~~

$$z=-\infty, g(z)=0$$

Decision boundary:

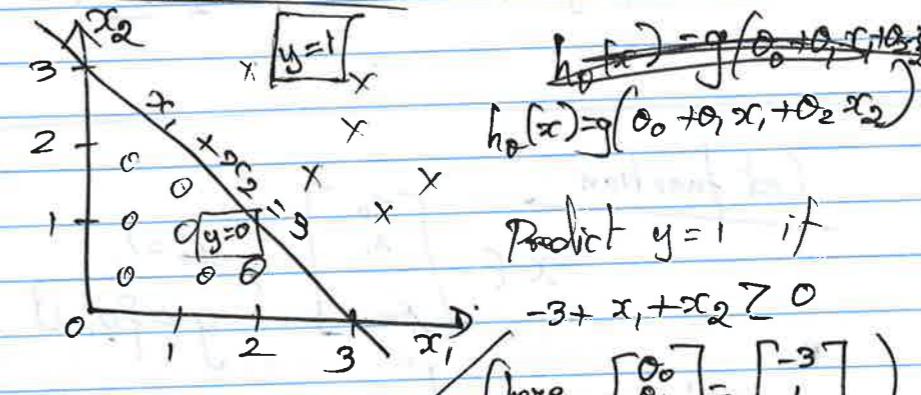
Say, $y=1$ if $h_\theta(x) \geq 0.5$, else $y=0$.



$$h_\theta(x) = g(\theta^T x) \rightarrow g(z) \geq 0.5$$

when $z \geq 0$
 $\Rightarrow h_\theta(x) \geq 0.5 (y=1)$
 when $\theta^T x \geq 0$.

And, $y=0$ when $\theta^T x < 0$.



$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

Predict $y=1$ if

$$-3 + x_1 + x_2 \geq 0$$

(here $\begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$)

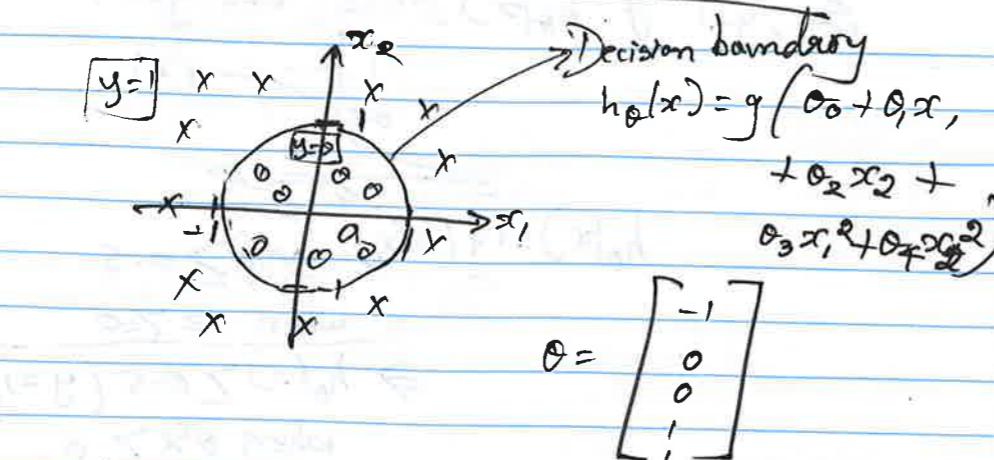
Since $y=1$ if $g(z) \geq 0.5$

$$\Rightarrow y=1 \text{ if } -3 + x_1 + x_2 \geq 0$$

$$\Rightarrow y=1 \text{ if } x_1 + x_2 \geq 3$$

Here, the line $x_1 + x_2 = 3$ is the decision boundary.

Non-linear decision boundary :



$$y=1 \text{ if } -1+x_1^2+x_2^2 \geq 0 \\ \Rightarrow y=1 \text{ if } x_1^2+x_2^2 \geq 1$$

Cost function :

$$x \in \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_{m-1} \end{bmatrix}, x_0 = 1, y \in \{0, 1\}$$

$$h_\theta(x) = \frac{1}{1+e^{-\theta^T x}}$$

From linear regression,

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(h_\theta(x), y)$$

$$\text{cost}(h_\theta(x), y) = \frac{1}{2} (h_\theta(x) - y)^2$$

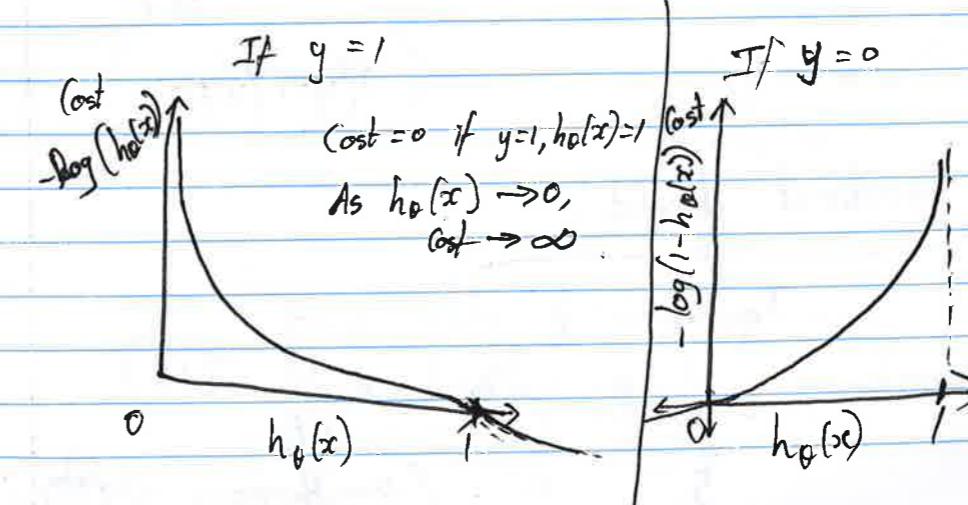
Cost function for logistic regression is
non-convex



Why is this cost function non-convex?

Logistic regression cost function :

$$\text{cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)), & \text{if } y=1 \\ -\log(1-h_\theta(x)), & \text{if } y=0 \end{cases}$$



Logistic regression Cost function (Simplified):

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y=1 \\ -\log(1-h_\theta(x)) & \text{if } y=0 \end{cases}$$

$$\Rightarrow \text{Cost}(h_\theta(x), y) = -y \log(h_\theta(x)) - (1-y) \log(1-h_\theta(x))$$

$$\therefore J(\theta) = \frac{1}{m} \sum_{i=1}^m$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})))$$

(Maximum Likelihood Estimation)

To fit params θ : $\min_{\theta} J(\theta)$

$$\text{To predict for a new } x: h_\theta(x) = \frac{1}{1+e^{-\theta^T x}} \quad (P(y=1/x; \theta))$$

Gradient descent:

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

{ (Simultaneous update for all)

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

$$\theta_j := \theta_j - \alpha \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

(Simultaneous update all θ_j)

Same as the gradient descent

for linear regression!!

But, $h_\theta(x)$ here is different.

Advanced Optimization:

Optimization techniques:

- Gradient descent

- Conjugate Gradient

- BFGS

- L-BFGS

No need to manually pick learning rate

Faster than gradient descent

More Complex

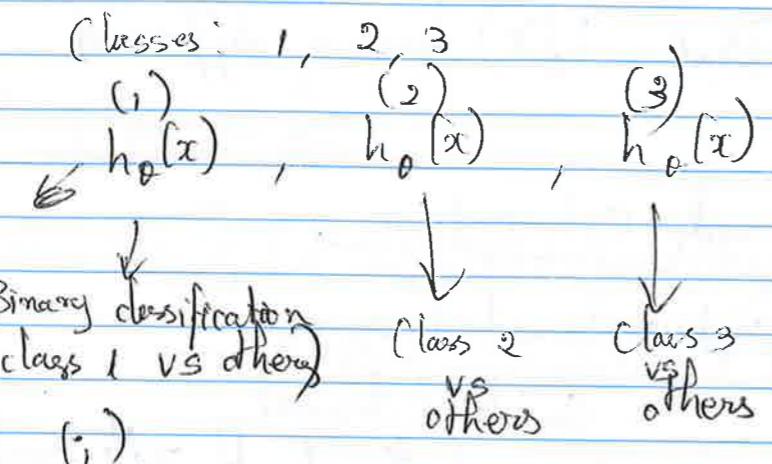
Essentially employs line search to automatically pick a good value of α (learning rate)

Logistic regression for multi-class classification:

→ one-vs-all (one-vs-rest) classification

(eg) Training set with three classes: A, B, C

→ Turn this into 3 separate binary classification problems



$h_0(x) \rightarrow$ predict $y=i$
 $(P(y=i/x; \theta))$

Output:

$$\max_i h_0^{(i)}(x)$$

Regularization:

- Underfitting (High bias)

- Overfitting (High variance)

"Just right"

To many hypothesis may fit the training set well (cost function is very low) but fails for new predictions.

Solutions for overfitting:

→ Reduce the no. of features

Manual selection / Model selection algorithm

→ Regularization

Keep all features

Reduce the magnitude of params θ_j

Works when there are a lot of features, each contributing a bit.

Penalize the cost function to make some θ_j values really small.

- Small values for parameters $\theta_0, \theta_1, \dots, \theta_n$
- "Simpler" hypothesis
 - Less prone to overfitting

~~We don't know which params to shrink,~~

(since we don't know which params correspond to which polynomial terms)

so regularize them all!

Regularized cost function (for linear regression):

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

$$\min_{\theta} J(\theta)$$

$\lambda \rightarrow$ Regularization parameter

θ_0 needsn't be regularized - since it's a constant term, not much impact.

What if λ is too large?

↳ Results in underfitting

(As $\theta_0, \theta_1, \dots, \theta_n$ will be almost zero)

Regularized linear regression:

$$\text{Cost fn., } J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

Gradient descent (with regularization):

$$\frac{\partial}{\partial \theta}$$

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_0^{(i)}$$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

$$\Rightarrow \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

$$\theta_j := \theta_j - \alpha \frac{1}{m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} + \lambda \cdot \theta_j \right]$$

$$\theta_j := \theta_j - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} - \alpha \frac{\lambda}{m} \theta_j$$

$$\theta_j := \theta_j \left(1 - \alpha \cdot \frac{\lambda}{m} \right) - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}$$

$$\left(1 - \alpha \cdot \frac{\lambda}{m} \right) < 1 \quad \text{Reduces } \theta_j \quad \text{Everything else is the same.}$$

Normal equation (with regularization)

$$X\theta = y$$

$$X^T X \theta = X^T y$$

$$\Rightarrow \theta = (X^T X)^{-1} X^T y$$

With regularization,

~~$\theta = (X^T X)^{-1}$~~

$$\theta = \left(X^T X + \lambda \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \right)^{-1} X^T y$$

$\hookrightarrow (n+1) \times (n+1)$ matrix

Even if $X^T X$ is non-invertible ($m \leq n$),

$$(X^T X + \lambda \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix})$$
 is always invertible

Regularized logistic regression:

$$\begin{aligned} J(\theta) = & - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) \right. \\ & \quad \left. + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \right] \\ & + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \end{aligned}$$

Gradient descent:

$$\theta_j := \theta_j \left(1 - \alpha \cdot \frac{\lambda}{m} \right) - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}))_j - y^{(i)} x_j^{(i)}$$

Week 4:

Neural networks: Representation:

Non-linear hypothesis:

- For training set with lot of features, polynomial feature mapping will generate too many features.

$n \rightarrow$ number of features

Number of quadratic features = $O(n^2)$

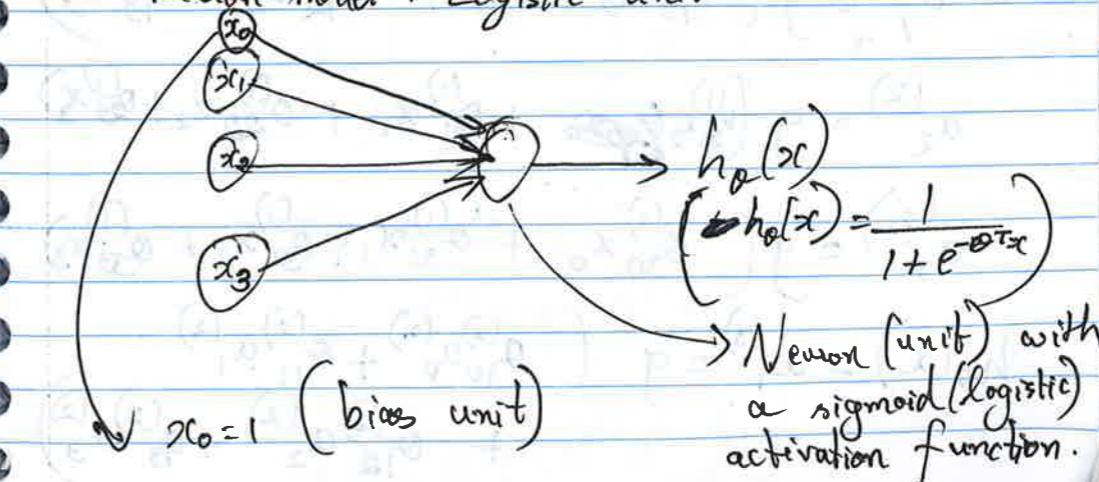
$$\text{(feature mapping)} \quad \left(\approx \frac{n^2}{2} \right)$$

$\hookrightarrow n^2$

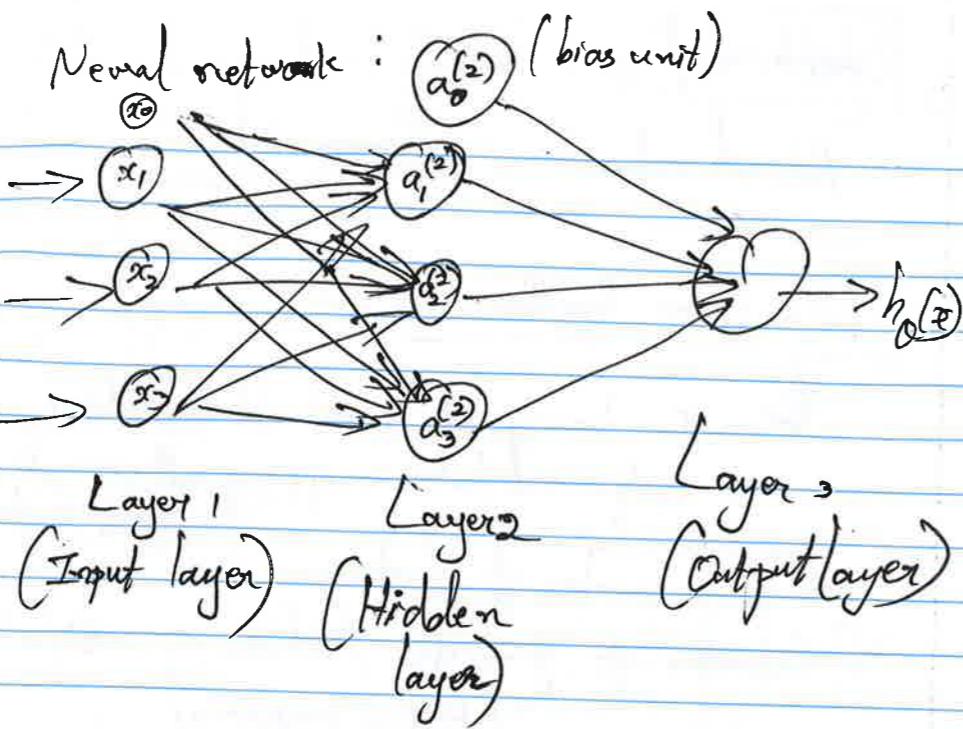
\hookrightarrow If number of features is large, total no. of generated features is very large even for quadratic hypothesis.

Model Representation:

Neuron model: Logistic unit



Neuron (unit) with a sigmoid (logistic) activation function.



Here, $\theta^{(1)} \in \mathbb{R}^{3 \times 4}$
 $\theta^{(2)} \in \mathbb{R}^{1 \times 4}$

Network Θ has s_j units in layer j
 and s_{j+1} units in layer $j+1$

s_j units in layer j
 &
 s_{j+1} units in layer $j+1$

$\Rightarrow \theta^{(3)}$ will be a
 $s_{j+1} \times (s_j + 1)$
 dimensional matrix

Mapping from
 layer j to layer $(j+1)$

Additional column for the
 bias unit.

$a_i^{(j)}$ \rightarrow activation of unit i
 in layer j

$\theta^{(j)}$ \rightarrow Matrix of weights
 controlling function mapping
 from layer j to layer $(j+1)$.

g \rightarrow sigmoid function (activation
 function)

$$a_1^{(2)} = g(\theta_{10}^{(1)}x_0 + \theta_{11}^{(1)}x_1 + \theta_{12}^{(1)}x_2 + \theta_{13}^{(1)}x_3)$$

$$a_2^{(2)} = g(\theta_{20}^{(1)}x_0 + \theta_{21}^{(1)}x_1 + \theta_{22}^{(1)}x_2 + \theta_{23}^{(1)}x_3)$$

$$a_3^{(2)} = g(\theta_{30}^{(1)}x_0 + \theta_{31}^{(1)}x_1 + \theta_{32}^{(1)}x_2 + \theta_{33}^{(1)}x_3)$$

$$h_0(x) = a_1^{(3)} = g(\theta_{10}^{(2)}a_0^{(2)} + \theta_{11}^{(2)}a_1^{(2)} + \theta_{12}^{(2)}a_2^{(2)} + \theta_{13}^{(2)}a_3^{(2)})$$

Forward prop: Vectorized implementation

$$a_1^{(2)} = g(z_1^{(2)}) \quad z_1^{(2)} = \theta_{10}^{(1)}x_0 + \theta_{11}^{(1)}x_1 + \theta_{12}^{(1)}x_2 \dots$$

$$a_2^{(2)} = g(z_2^{(2)})$$

$$a_3^{(2)} = g(z_3^{(2)})$$

$$z_1^{(2)} = \theta_{11}^{(1)}x$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}; \quad z^{(2)} = \theta^{(1)}x$$

$$\Rightarrow a^{(2)} = g(z^{(2)})$$

$$\begin{cases} a^{(2)} = g(\theta^{(1)}.x) \\ \text{(or)} \\ a^{(2)} = g(\theta^{(1)}.a^{(1)}) \end{cases} \quad (\alpha^1 = x).$$

Now, add the bias unit

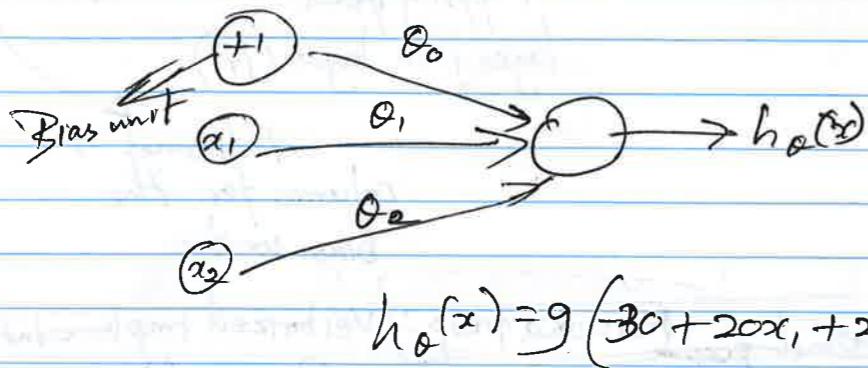
Add $a_0^{(2)} = 1$.

$$z^{(3)} = \theta^{(2)} \cdot a^{(2)}$$

$$h_{\theta}(x) = a^{(3)} = g(z^{(3)}) = g(\theta^{(2)} \cdot a^{(2)})$$

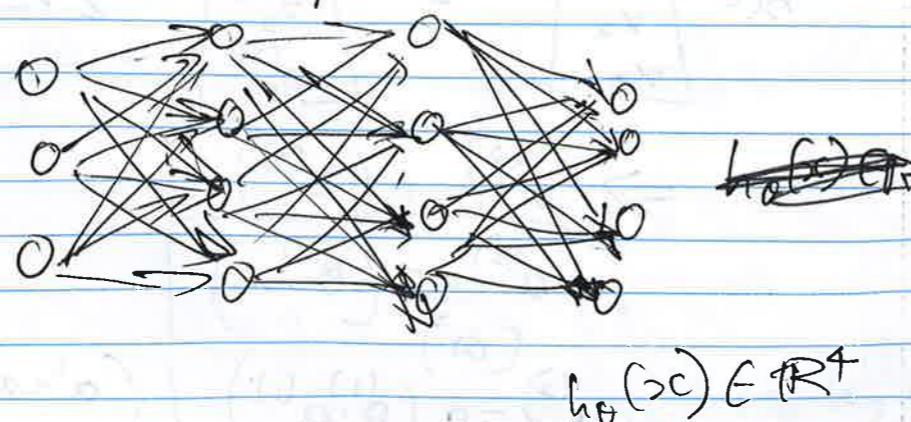
Example: AND

$$x_1, x_2 \in \{0, 1\} ; y = x_1 \text{ AND } x_2$$



Multi-class classification neural network:

- One vs all
- multiple output units



Training set: $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots$

$$y^{(i)} \text{ one of } \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

Class 1 Class 2 Class 3 Class 4

Week 5

Neural n/w (classification):

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

↳ Training set

L = total no. of layers in NN

s_l = no. of units (not including the bias unit) in layer l .

Cost function:

Logistic regression cost fn.:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2.$$

Cost function of NN:

$K \geq$ no. of classes (i.e., no. of output units)

$$\begin{aligned} h_{\theta}(x) &\in \mathbb{R}^k \\ (h_{\theta}(x))_i &= i^{th} \text{ output} \end{aligned}$$

$$\begin{aligned} J(\theta) &= -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\theta}(x^{(i)}))_k \right. \\ &\quad \left. + (1-y_k^{(i)}) \log(1-h_{\theta}(x^{(i)}))_k \right] \\ &\quad + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^m \sum_{j=1}^{s_l} (\theta_{ji}^{(l)})^2 \end{aligned}$$

Backpropagation:

$$\min_{\theta} J(\theta)$$

Need code to compute:

$$J(\theta) \text{ & } \frac{\partial}{\partial \theta_{ji}} J(\theta)$$

Gradient computation: Backpropagation algorithm:

Intuition: $\delta_j^{(l)}$ \rightarrow error of ~~unit~~ unit j in layer l .
(i.e.) error in $a_j^{(l)}$.

For each output unit (layer L)

$$\delta_j^L = a_j^L - y_j$$

$$(\rightarrow L \rightarrow (h_{\theta}(x))_j)$$

$$\rightarrow \text{Vectorized: } \delta^L = a^L - y$$

↓ Errors in the output layer.

Now compute the errors in the previous layers.

$$\delta^{(L-1)} = ((\theta^{(L-1)})^T \cdot \delta^{(L)}) \otimes g'(z^{(L-1)})$$

Element-wise multiplication

NN with 4 layers:

$$\begin{aligned} \delta^{(4)} &= a^{(4)} - y \\ \delta^{(3)} &= (\theta^{(3)})^T \delta^{(4)} \otimes g'(z^{(3)}) \\ \delta^{(2)} &= (\theta^{(2)})^T \delta^{(3)} \otimes g'(z^{(2)}) \end{aligned}$$

$$\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta) = a_j^{(l)} \delta_i^{(l+1)}$$

(Can be proven)

Ignoring λ here
($\lambda=0$)
No regularization

$$\rightarrow D_{ij}^{(l)} := \frac{1}{m} (\Delta_{ij}^{(l)} + \lambda \theta_{ij}^{(l)}) \text{ if } j \neq 0$$

$$\rightarrow D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \quad (\text{No regularization for the bias unit}).$$

$$\frac{\partial}{\partial \theta_{0j}^{(l)}} J(\theta) = D_{0j}^{(l)}$$

Backpropagation algorithm:

Training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

Set $\Delta_{ij}^{(l)} = 0$ (for all i, j, l)
(Used to compute
 $\frac{\partial}{\partial \theta_{ij}^{(l)}} J(\theta)$)

For $i = 1$ to m (each training set $(x_i^{(i)}, y^{(i)})$):

→ Set $a^{(1)} = x^{(i)}$

→ Perform forward prop to compute

$a^{(l)}$ for $l = 2, 3, \dots, L$.

→ Using $y^{(i)}$, compute $\delta^{(L)} = a^{(L)} - y^{(i)}$

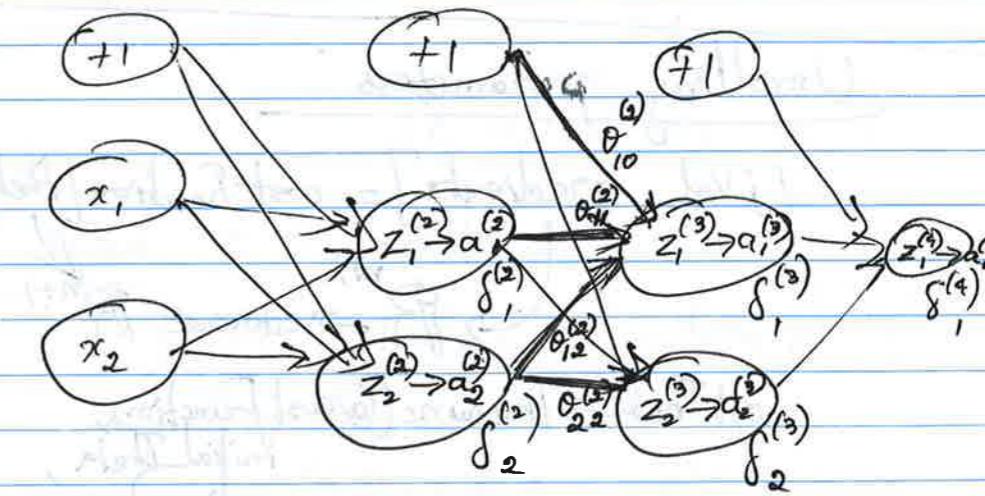
→ Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

→ $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

↳ Vectorized: $\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (a^{(l)})^T$

Back prop (Intuition):

forward prop:



$$z_1^{(3)} = \theta_{10}^{(2)} x_1 + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)}$$

Cost function (for 1 training set, with one-class (one output unit), i.e., $k=1$, and no regularization, i.e., $\lambda=0$)

$$\text{cost}(i) = y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)}))$$

$\delta_j^{(l)}$ → error of cost for $a_j^{(l)}$ (unit j in layer l).
 $\delta_j^{(l)} = \frac{\partial}{\partial z_j^{(l)}} (\text{cost}(i))$ (for $j >= 0$)

$$\delta_i^{(4)} = y_i - a_i^{(4)}$$

$$\delta_2^{(2)} = \theta_{12}^{(2)} \cdot \delta_1^{(3)} + \theta_{22}^{(2)} \cdot \delta_2^{(3)}$$

$$\delta_2^{(3)} = \theta_{12}^{(3)} \cdot \delta_1^{(4)}$$

Unrolling parameters:

[j Val, gradients] = costFunction(theta)

→ $\theta \in \mathbb{R}^n$ → vectors $\in \mathbb{R}^{n+1}$

optTheta = fminunc(@(costFunction, initialTheta, options)

NN (4 layers):

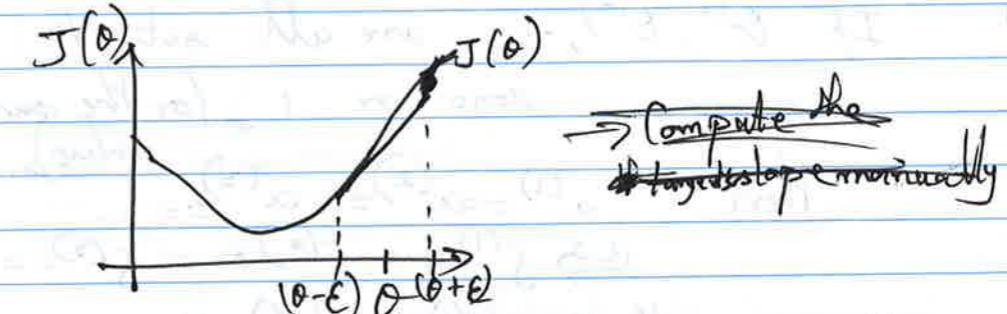
$\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$ - matrices of params

$D^{(1)}, D^{(2)}, D^{(3)}$ - gradient matrices

Unroll them into
(column-major) Vectors for θ

Gradient Checking:

- Numerical estimation of gradients
- To verify back prop works as expected



Manually compute the slope of the tangent

$$\frac{d J(\theta)}{d \theta} \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

(ϵ is typically 10^{-4})

When θ is a vector:

$\theta \in \mathbb{R}^n$ (after unrolling $\theta^{(1)}, \theta^{(2)}, \dots$)

$$\theta = [\theta_1, \theta_2, \dots, \theta_n]$$

$$\frac{\partial}{\partial \theta_1} J(\theta) \approx \frac{(J(\theta_1 + \epsilon, \theta_2, \theta_3, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \theta_3, \dots, \theta_n))}{2\epsilon}$$

$$\frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{(J(\theta_1, \theta_2 + \epsilon, \theta_3, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \theta_3, \dots, \theta_n))}{2\epsilon}$$

Random Initialization :

Initial value of θ

zero initialization doesn't work.

If $\theta^{(1)}, \theta^{(2)}, \dots$ are all set to zero or 1 (or the same value),
 then $a^{(1)} = a^{(2)} = a^{(3)} = \dots$
 $\hookrightarrow f^{(1)} = f^{(2)} = f^{(3)} = \dots$
 $\hookrightarrow J(\theta) = J(\theta_2) = \dots$
 $\hookrightarrow \frac{\partial J(\theta)}{\partial \theta_1^{(1)}} = \frac{\partial J(\theta)}{\partial \theta_2^{(1)}} = \dots$
 \hookrightarrow After 1 iteration $\theta^{(2)} = \theta^{(3)} = \dots$

Use random initialization

Symmetry breaking

Initialize each θ_{ij}^l to random value
in $[-\epsilon, \epsilon]$

Steps:

- Pick NN architecture/weights
- Randomly initialize weights
- Compute $h_\theta(x)$ (Forward prop)
- Compute cost function $J(\theta)$
- Compute gradients $\frac{\partial}{\partial \theta_k} J(\theta)$

-
- Compute backprop (Back prop)
gradient checking and verify it works
- Minimize $J(\theta)$ using gradient descent or advanced optimization

NN cost function need not be convex

→ It's okay to converge to local minima.

Week 6

Evaluating a hypothesis :

- Overfitting

→ Split dataset into training set & test set

Misclassification error (0/1 misclassification err)

$$\text{err}(h_\theta(x), y) = \begin{cases} 1, & \text{if } h_\theta(x) \geq 0.5, y=0 \\ & \text{or if } h_\theta(x) < 0.5, y=1 \\ 0, & \text{otherwise} \end{cases}$$

$$\text{Test error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_\theta(x_{\text{test}}^{(i)}), y_{\text{test}}^{(i)})$$

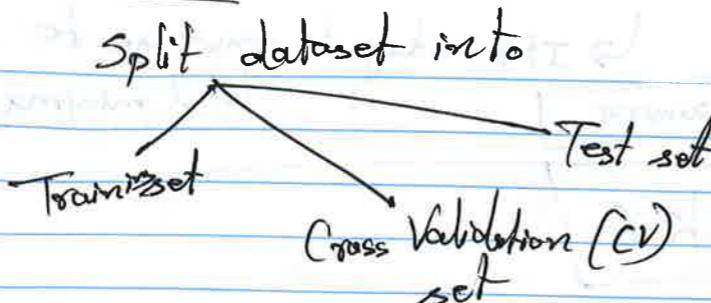
Model Selection:

- Picking a function for $h_\theta(x)$
- What degree of polynomial to pick?

Say, pick a dg degree (d) based on which minimizes $J_{\text{test}}(\theta)$.

Not fair as $J_{\text{test}}(\theta)$ is likely to be an optimistic estimate of generalization error (i.e., d is chosen to fit the test set)

Solution:



- Branch of models
- Pick the best ~~best~~ performing model (lowest ~~cost~~ function or error) on the cross validation set (i.e., pick the model with lowest cross validation error)

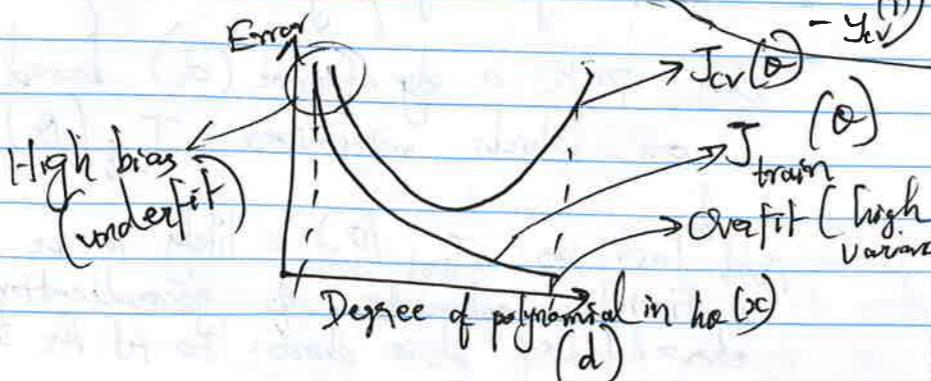
Bias and Variance:

High Bias - underfitting

High Variance - overfitting

$$\text{Training error: } J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\text{(cross-validation error: } J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2)$$



Underfit (bias):

$$J_{\text{train}}(\theta) \text{ will be high}$$

$$J_{\text{cv}}(\theta) \approx J_{\text{train}}(\theta)$$

Overfit (variance):

$$J_{\text{train}}(\theta) \text{ will be low}$$

$$J_{\text{cv}}(\theta) \gg J_{\text{train}}(\theta)$$

Regularization and bias/variance:

Large λ (regularization param)

→ Underfit

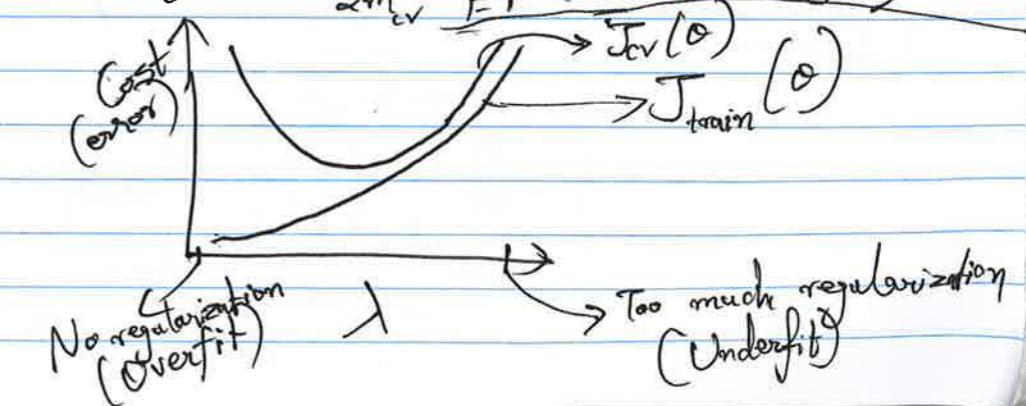
Small $\lambda \rightarrow$ Overfit

Choosing λ correctly:

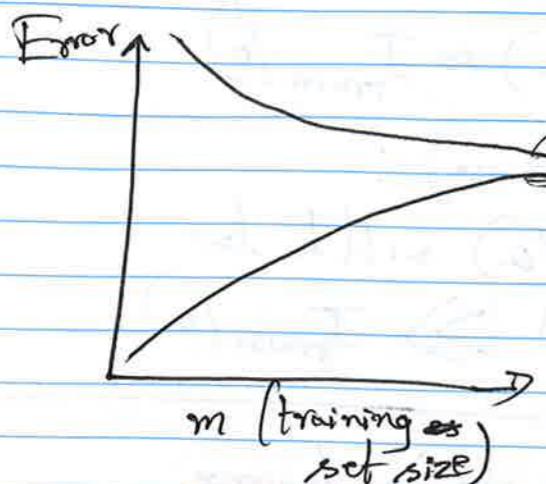
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2$$

$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

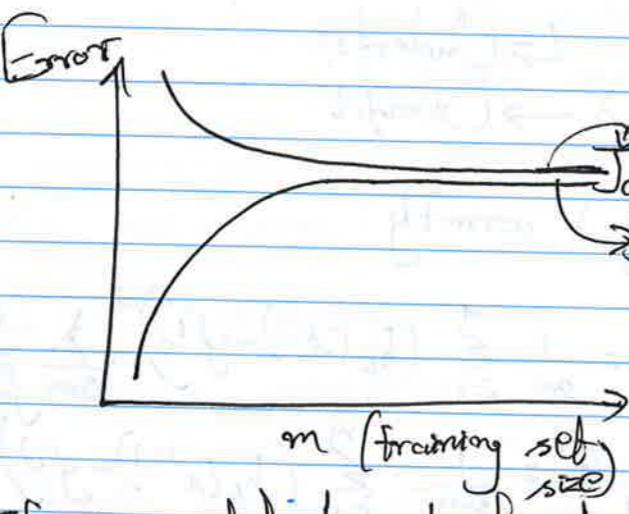
$$J_{\text{cv}}(\theta) = \frac{1}{2m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x_{\text{cv}}^{(i)}) - y_{\text{cv}}^{(i)})^2$$



Learning Curves:

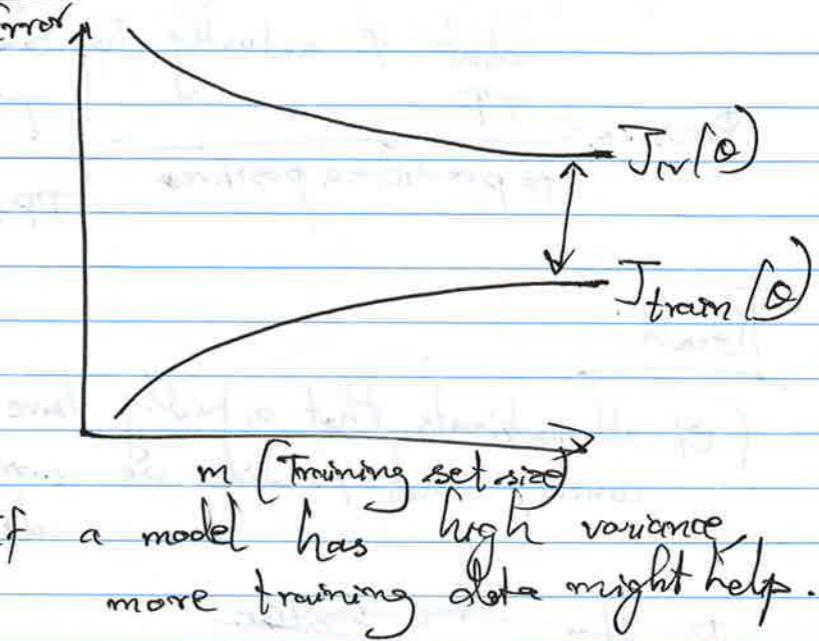


High bias (Underfit):



If a model has high bias,
more training data won't help.

High variance (Overfit):



If a model has high variance,
more training data might help.

Error metrics for skewed classes:

Skewed class \rightarrow Training set has skew in output labels

99.5% $y=0$, 0.5% $y=1$

Precision / Recall:

$y=1$ in presence of rare class that we want to predict.

		Actual class	
		1	0
Predicted class	1	TP	FP
	0	FN	TN

Precision:

(Of all patients where we predicted $y=1$, what % actually has cancer)

$$\text{Precision} = \frac{TP}{\# \text{predicted positives}} = \frac{TP}{TP+FP}$$

Recall:

(Of all patients that actually have cancer, what % did we correctly detect)

Recall = True positive

$$\text{Recall} = \frac{TP}{\# \text{actual positives}} = \frac{TP}{TP+FN}$$

Trading off precision and recall:

→ Say, predict cancer only if highly confident.

Predict 1 if $h_{\theta}(x) >= 0.7$

Predict 0 if $h_{\theta}(x) < 0.7$

↳ Higher precision, lower recall.

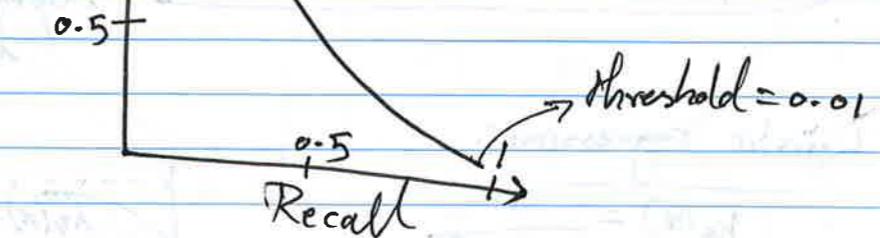
→ Say, we want to avoid missing too many cases of cancer (avoid false negatives)

Predict 1 if $h_{\theta}(x) >= 0.3$, 0 otherwise

↳ Lower precision, higher recall

Predict 1 if $h_{\theta}(x) >= \text{threshold}$

Precision
1 → threshold = 0.99



F_1 score (F score):

- Comparing precision/recall numbers

	Precision (P)	Recall (R)
Algo 1	0.5	0.4
Algo 2	0.7	0.1
Algo 3	0.02	1.0

$$F_1 \text{ score} = \frac{2PR}{P+R}$$

$P=0$ or $R=0 \Rightarrow F\text{-score}=0$

$P=1$ and $R=1 \Rightarrow F\text{-score}=1$

Week 7

SVM (Support Vector Machine):

Optimization Objective:

Supervised learning

Logistic regression:

$$h_0(x) = \frac{1}{1+e^{-\theta^T x}} ; \quad h_0(x) = g(z)$$

If $y=1$, we want $h_0(x) = 1$, $\theta^T x \gg 0$

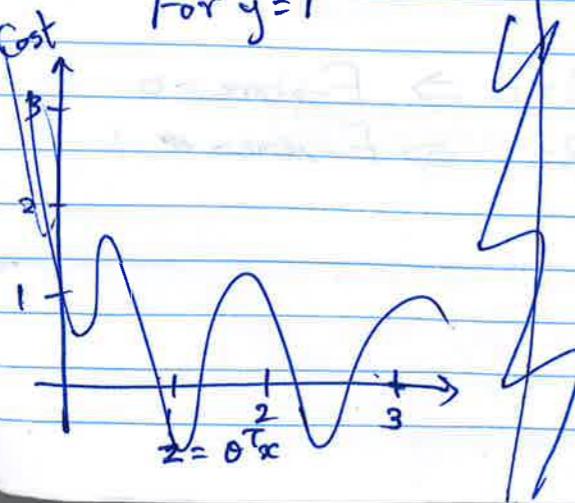
If $y=0$, we want $h_0(x) \approx 0$, $\theta^T x \ll 0$

$$\text{Cost} = -y \log h_0(x) - (1-y) \log(1-h_0(x))$$

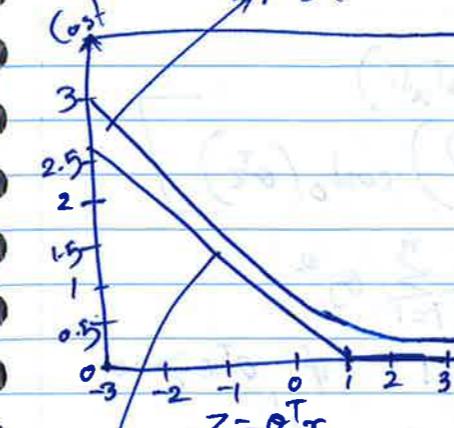
$$\text{Cost} = \left(-y \log \left(\frac{1}{1+e^{-\theta^T x}} \right) \right) + \left((1-y) \log \left(1 - \frac{1}{1+e^{-\theta^T x}} \right) \right)$$

For $y=1$

For $y=0$

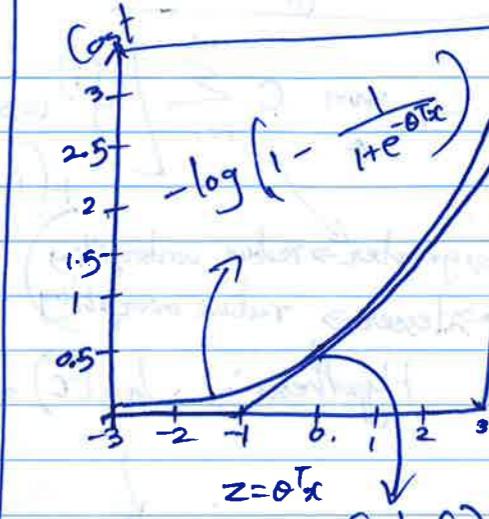


$$y=1 : \log \left(\frac{1}{1+e^{-\theta^T x}} \right)$$



$$\text{cost}_1(z) = \max(0, k(1-z))$$

$$y=0 :$$



$$\text{cost}_0(z) = \max(0, k(1+z))$$

SVM cost function:

$$\min_{\theta} \sum_{i=1}^m \left(y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right) + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

\Leftrightarrow No division by m
 → Regularization parameter moved
 (Instead of $A + \lambda B$,
 use $CA + \lambda B$.
 $(C = \frac{1}{\lambda})$)

SVM hypothesis:

$$\rightarrow \min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right]$$

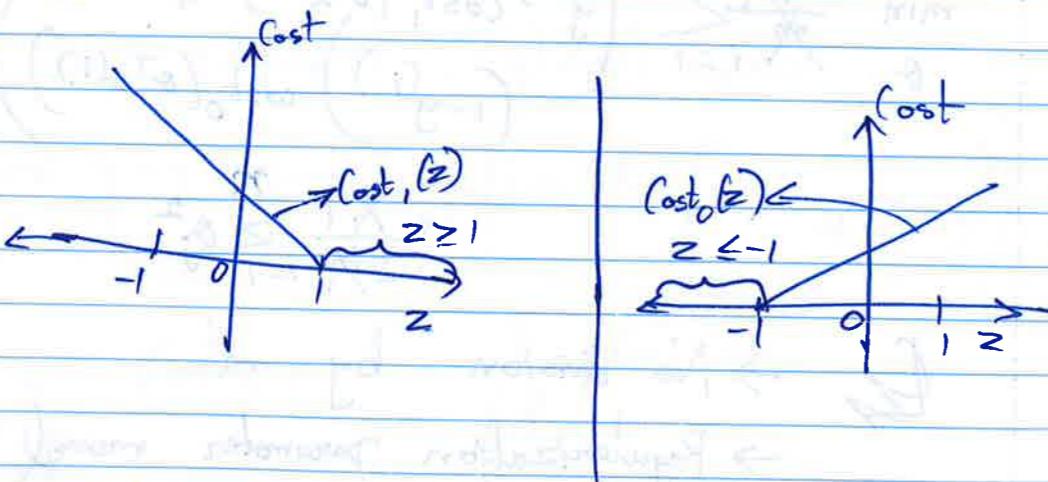
($C \rightarrow$ greater \rightarrow reduce underfitting)
($C \rightarrow$ lesser \rightarrow reduce overfitting) $+ \frac{1}{2} \sum_{j=1}^n \theta_j^2$

Hypothesis: $h(\theta)(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$

Does not predict the probability

(Discriminant function) like logistic regression

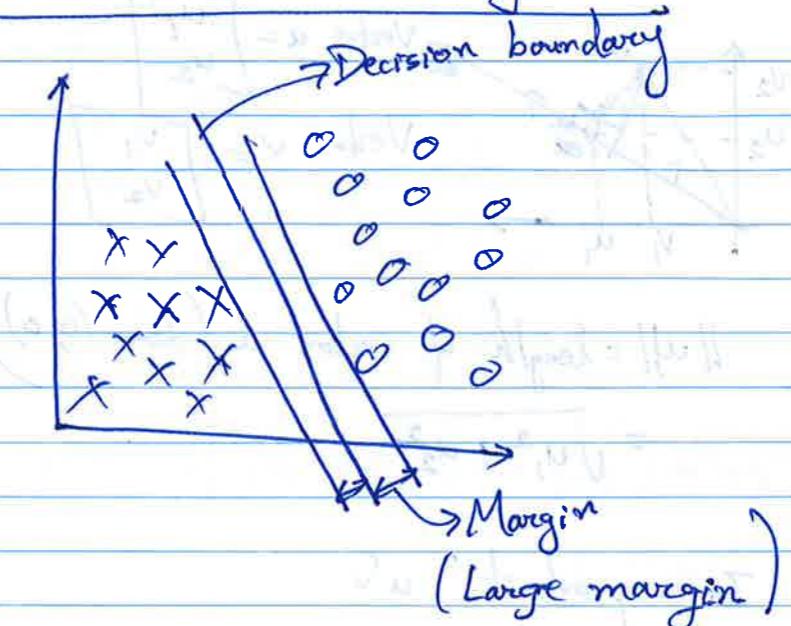
SVM \rightarrow Large Margin Classifier:



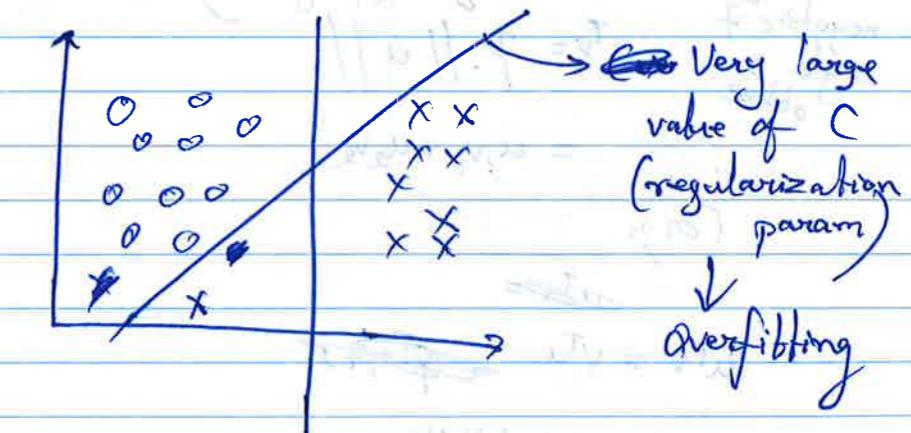
If $y=1$, we want $\theta^T x \geq 1$ (not just $\theta^T x \geq 0$)

If $y=0$, we want $\theta^T x \leq -1$ (not just $\theta^T x \leq 0$)

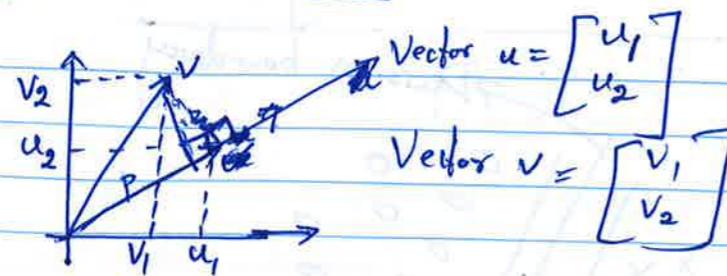
SVM Decision Boundary:



Large margin classifier in presence of outliers:



Vector Inner product:



$\|u\| = \text{length of vector } u \text{ (from } (0,0))$

$$= \sqrt{u_1^2 + u_2^2}$$

Inner product $u^T v$

$P \rightarrow \text{length of projection of } v \text{ onto } u.$
 negative if angle is obtuse
 $u^T v = P \cdot \|u\|$
 can be signed (negative)

$$= u_1 v_1 + u_2 v_2$$

(or)

~~$u^T v = v^T u$~~

$$= v_1 u_1 + v_2 u_2$$

SVM Decision boundary:

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

$$\text{s.t., } \theta^T x^{(i)} \geq 1 \text{ if } y^{(i)} = 1$$

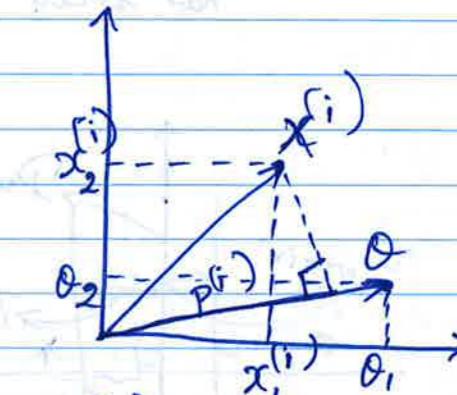
$$\theta^T x^{(i)} \leq -1 \text{ if } y^{(i)} = 0$$

Simplification: $\theta_0 = 0, n=2$

$$\Rightarrow \min_{\theta} \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} \left(\sqrt{\theta_1^2 + \theta_2^2} \right)^2 = \frac{1}{2} \|\theta\|^2$$

$$\theta^T x^{(i)} = ?$$

$\uparrow \quad \uparrow$
 $u^T \quad v$



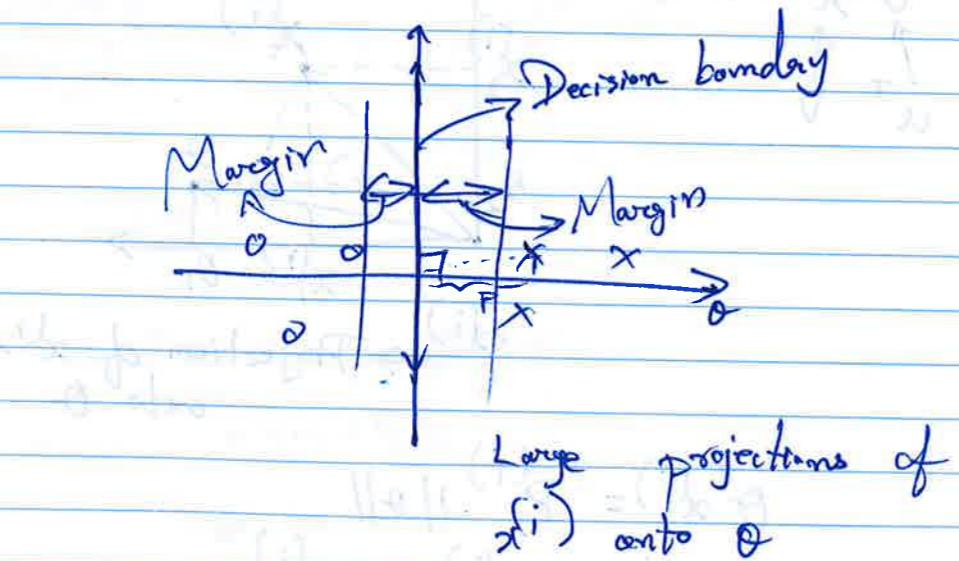
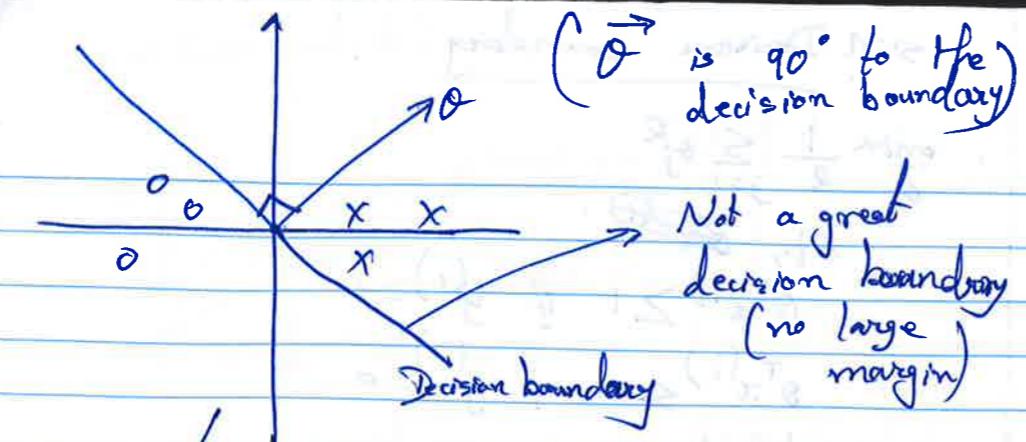
$P^{(i)} \rightarrow \text{Projection of } x^{(i)} \text{ onto } \theta$

$$\theta^T x^{(i)} = P^{(i)} \cdot \|\theta\|$$

$$= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$$

$$\Rightarrow \min_{\theta} \frac{1}{2} \|\theta\|^2 \text{ s.t., } P^{(i)} \cdot \|\theta\| \geq 1 \text{ if } y^{(i)} = 1,$$

$$P^{(i)} \cdot \|\theta\| \leq -1 \text{ if } y^{(i)} = 0$$

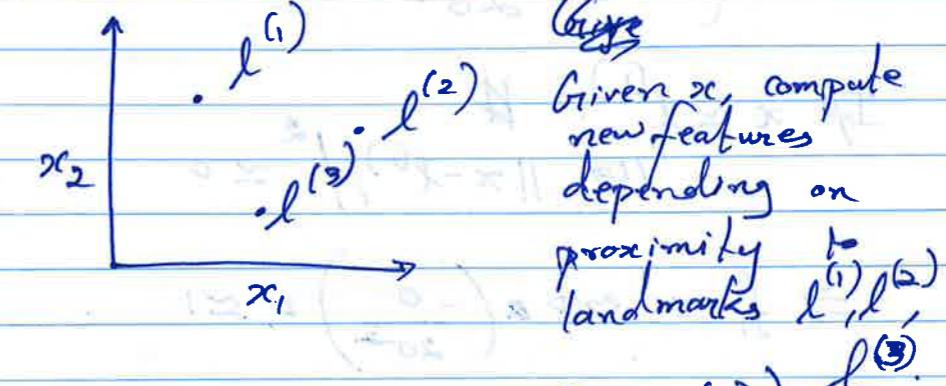


Non-linear decision boundary:

Option 1: Polynomial features.
 or Kernel:

(Predict $y = 1$ if
 $\theta_0 x_1 + \theta_1 x_1^2 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_2^2 + \theta_5 x_1 x_2 + \dots \geq 0$)

Kernel:



Given x : $f_1 = \text{similarity}(x, l^{(1)})$

$$= \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

$f_2 = \text{similarity}(x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$

$f_3 = \text{similarity}(x, l^{(3)}) \dots$

kernel (this similarity is a Gaussian Kernel)

Kernels and Similarity:

$$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{20^2}\right)$$
$$= \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^{(1)})^2}{20^2}\right)$$

If $x \approx l^{(1)}$: $\Downarrow \|x - l^{(1)}\|^2 \approx 0$

$$\Rightarrow f_1 \approx \exp\left(-\frac{0}{20^2}\right) \approx 1$$

If x is far from $l^{(1)}$:

$$f_1 \approx \exp\left(-\frac{(\text{large number})^2}{20^2}\right) \approx 0$$

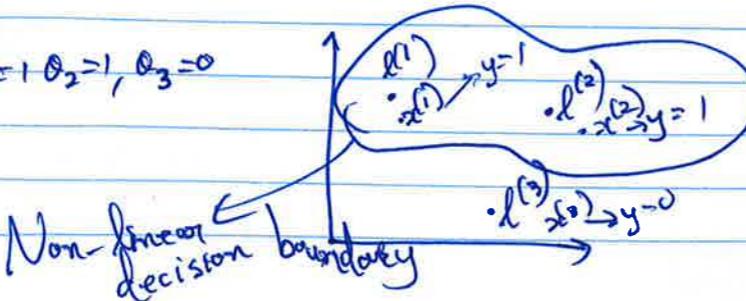
New features $\rightarrow f_1, f_2, f_3$

Predict $y=1$ when

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$

Example

$$\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$$



Choosing the landmarks:

How do you choose $l^{(1)}, l^{(2)}, l^{(3)}$?

SVM with kernels:

Given $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$,

choose $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$.

Given example x :

$f_1 = \text{similarity}(x, l^{(1)})$

$f_2 = \text{similarity}(x, l^{(2)})$

$$f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix}, \quad f_0 = 1$$

Feature vector

For training example $(x^{(i)}, y^{(i)})$:

$$f_1^{(i)} = \text{similarity}(x^{(i)}, l^{(1)})$$

$$f_2^{(i)} = \text{similarity}(x^{(i)}, l^{(2)})$$

$$\vdots$$
$$f_m^{(i)} = \text{similarity}(x^{(i)}, l^{(m)}) = \exp\left(-\frac{0}{20^2}\right) = 1$$

$$f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix}$$

$$f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix}; f_0^{(i)} = 1$$

Feature vector for $x^{(i)}$

Hypothesis: Given x^c , compute features $f \in \mathbb{R}^{m+1}$.
Predictor: $y=1$ if $\theta^T f \geq 0$

$$\theta^T f = \theta_0 f_0 + \theta_1 f_1 + \dots + \theta_m f_m$$

$\theta \in \mathbb{R}^{m+1}$

Training:

$$\min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \text{cost}_+(\theta^T f^{(i)}) + (1 - y^{(i)}) \text{cost}_-(\theta^T f^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

(Note, no. of features $f =$
no. of training examples
(+1 for bias)).

SVM parameters:

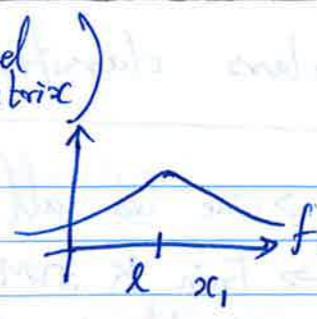
$C (= \frac{1}{\lambda}) \rightarrow$ large $C \rightarrow$ low bias, high variance (overfit)
 $C \rightarrow$ small $C \rightarrow$ high bias, low variance (underfit)

σ^2 (in Gaussian kernel for similarity matrix)

σ^2 (for Gaussian kernel in similarity matrix)

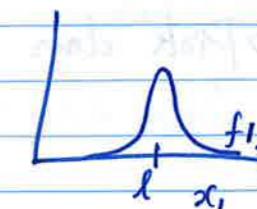
Large σ^2 :

Underfitting



Small σ^2 :

Overfitting



Using an SVM:

→ Choose param C

→ Choose kernel (similarity function)

→ No kernel (linear kernel):

Predict $y=1$ if $\theta^T x \geq 0$

→ Gaussian kernel:

$$f_i = \exp\left(-\frac{\|x - x^{(i)}\|^2}{2\sigma^2}\right), \text{ where } x^{(i)} = x^{(i)}$$

(Need to choose σ^2).

→ Perform feature scaling before using the kernel.

→ Polynomial kernels: $k(x, l) = (x^T l)^2, (x^T l)^3, (x^T l + 1)^3, (x^T l + 5)^4$

Multiclass classification with SVM:

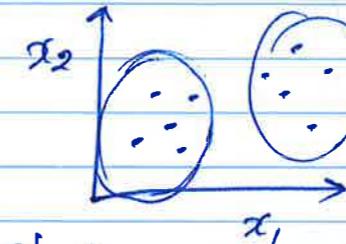
- One vs all method
- Train K SVMs, each for one of the K classes.
- Pick class i with largest $(\alpha^{(i)})^T x_c$.

Logistic Regression vs ~~SVM~~ SVM:

- $n \rightarrow$ no. of features
- $m \rightarrow$ no. of training examples
- large n (relative to m):
 - ↳ Use logistic regression or SVM w/o kernel (linear kernel)
- small n , intermediate m :
 - ↳ SVM w/ Gaussian kernel
- small n , large m :
 - ↳ Add more features, then use logistic regression or SVM w/o kernel.

Week 8:

Unsupervised Learning:



Training set:
 $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

- Clustering No labels of i !

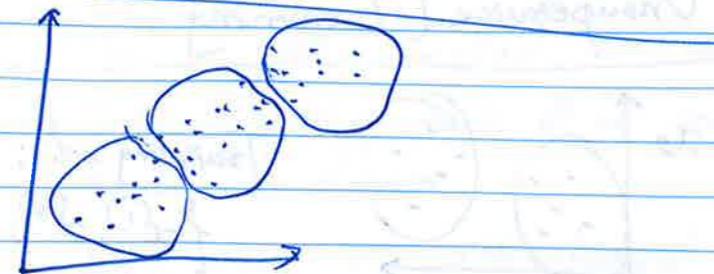
Clustering: K-means algorithm:

- Input:
- K (number of clusters)
 - Training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$
 - $x^{(i)} \in \mathbb{R}^n$ (drop $x_0=1$ convention)

Algo:

- Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$
- Repeat {
 - for $i = 1$ to m
 - $c^{(i)} :=$ index (from 1 to K) of cluster centroid closest to $x^{(i)}$
- Cluster assignment step } for $k = 1$ to K
 - $\mu_k =$ average (mean) of points assigned to cluster k .
- Move centroid step }

K-means for non-separated clusters :



Optimization Objective :

$c^{(i)}$ → cluster for example $x^{(i)}$
 μ_k → centroid of cluster k .
 $\mu_{c(i)}$ → centroid of cluster to which
 $x^{(i)}$ has been assigned

Cost fn. :

$$J(c^{(1)}, c^{(2)}, \dots, c^{(m)}, \mu_1, \mu_2, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c(i)}\|^2$$

$$\min_{c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

Random Initialization :

- Random initialization of cluster centroids at startup
- Should have $K < m$
- Randomly pick K training examples
- Set $\mu_1, \mu_2, \dots, \mu_K$ equal to these K examples.

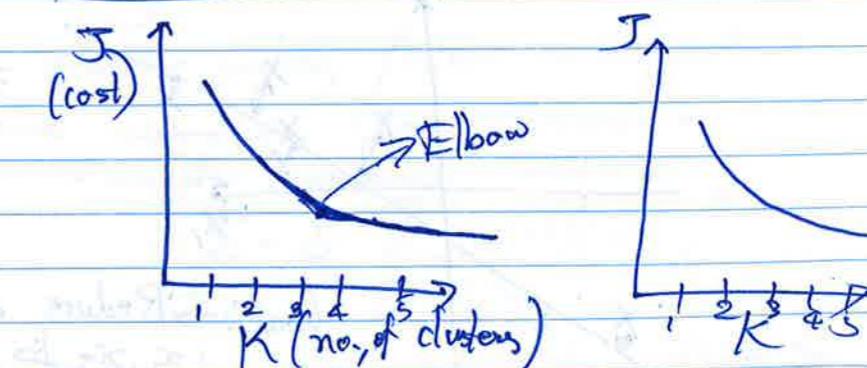
Note:

Local optima is possible based on centroid initialization.

Solution : → Try K-means multiple times with different random centroid initialization.
 → Pick the one with the least cost.

Choosing K (no. of clusters) :

- Manually, visualization, etc.
- Elbow method :

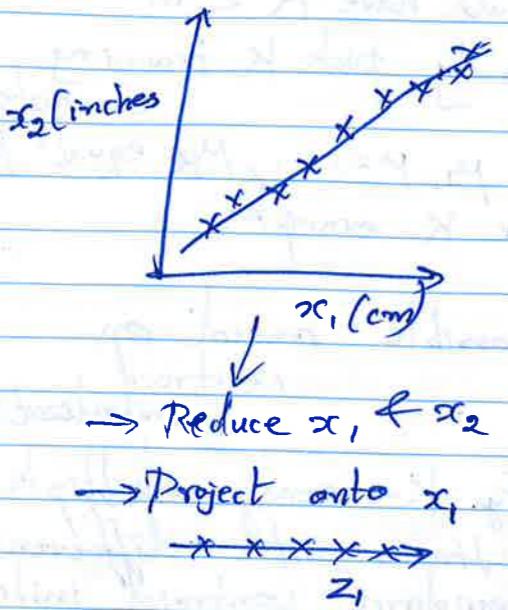


Dimensionality Reduction:

~~Method~~ ↗

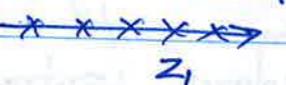
- Data Compression:

→ Reduce unnecessary or redundant features



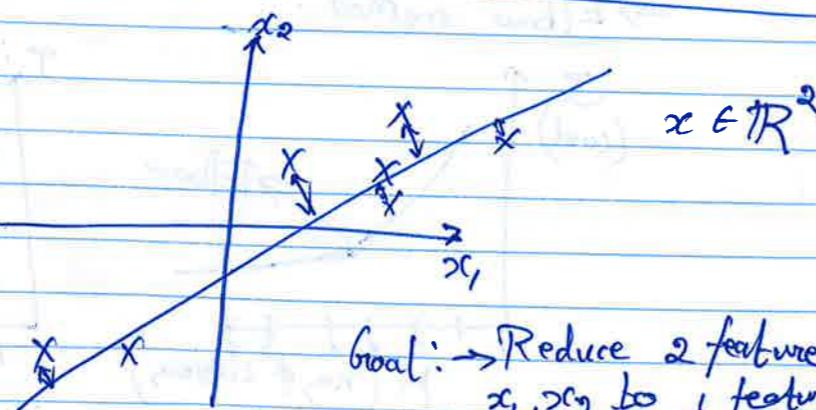
→ Reduce x_1 & x_2 into a new feature z_1

→ Project onto x_1



- Data Visualization:

PCA (Principal Component Analysis):



Goal: → Reduce 2 feature x_1, x_2 to 1 feature

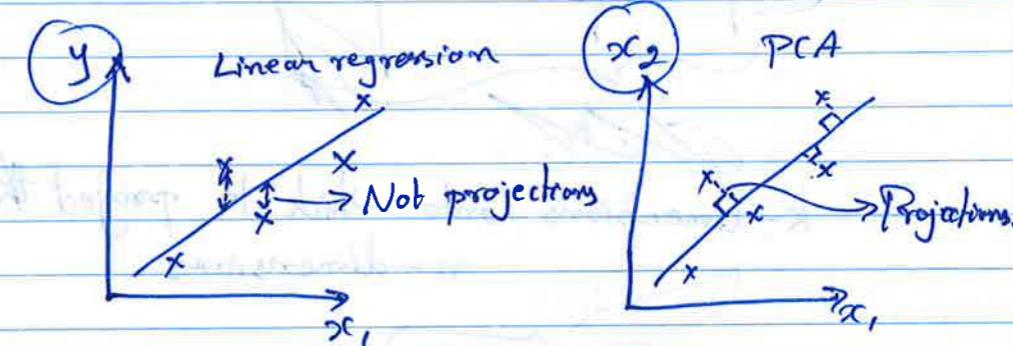
→ Find a line to project the points onto

PCA → find a ~~line~~ line such that the projection error (~~that~~ sum of squares of projection distance) is minimized.

Reduce n-dimensions to k-dimensions:

→ Find k vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project the data, so as to minimize the projection error.

PCA is NOT linear regression:



PCA algorithm: → How to find the k dimensions
Then, how to find the projected new features for each training example

→ Data preprocessing:

Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$

→ Feature scaling / Mean normalization

→ Reduce n-dimensions to k-dimensions

→ Compute covariance matrix Σ

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)})^T (x^{(i)})$$

(covariance matrix)

→ Compute eigen vectors of matrix Σ (Singular Value Decomposition)

$$[U, S, V] = svd(\Sigma)$$

Covariance matrix Σ
 \downarrow
 $n \times n$ matrix.

Covariance \rightarrow How much two random variables change together.

$$[U, S, V] = \text{svd}(\Sigma)$$

$n \times n$

$$V = \begin{bmatrix} u^{(1)} & u^{(2)} & \dots & u^{(n)} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

k -dimensions onto which to project the n -dimensions.

$$V = \begin{bmatrix} | & | & | & | \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \\ | & | & | & | \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$$x \in \mathbb{R}^n$$

$$x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^k$$

$$\begin{bmatrix} u_{(1)} & u_{(2)} & \dots & u_{(k)} \\ | & | & \dots & | \end{bmatrix} \in \mathbb{R}^{n \times k}$$

V_{reduce}

$$z^{(i)} = V_{\text{reduce}} \cdot x^{(i)}$$

$(k \times n) \cdot (n \times 1) \Rightarrow k \times 1 \Rightarrow k$ -dimensions.

$$z^{(i)} \in \mathbb{R}^k$$

Summary of PCA algo:

$$1.) X \in \mathbb{R}^{m \times n}$$

\rightarrow covariance matrix

$$2.) \Sigma = \frac{1}{m} \cdot X \cdot X^T$$

$$3.) [U, S, V] = \text{svd}(\Sigma)$$

$$4.) V_{\text{reduce}} = V[:, 1:k]$$

First k columns of V

$$5.) z = X \cdot V_{\text{reduce}}$$

$$\mathbb{R}^{m \times k} \quad \mathbb{R}^{n \times k}$$

Reconstruction from Compressed Representation:

Given $z \in \mathbb{R}^k$, get back $x \in \mathbb{R}^n$.

$(\text{approx.}, \approx x)$.

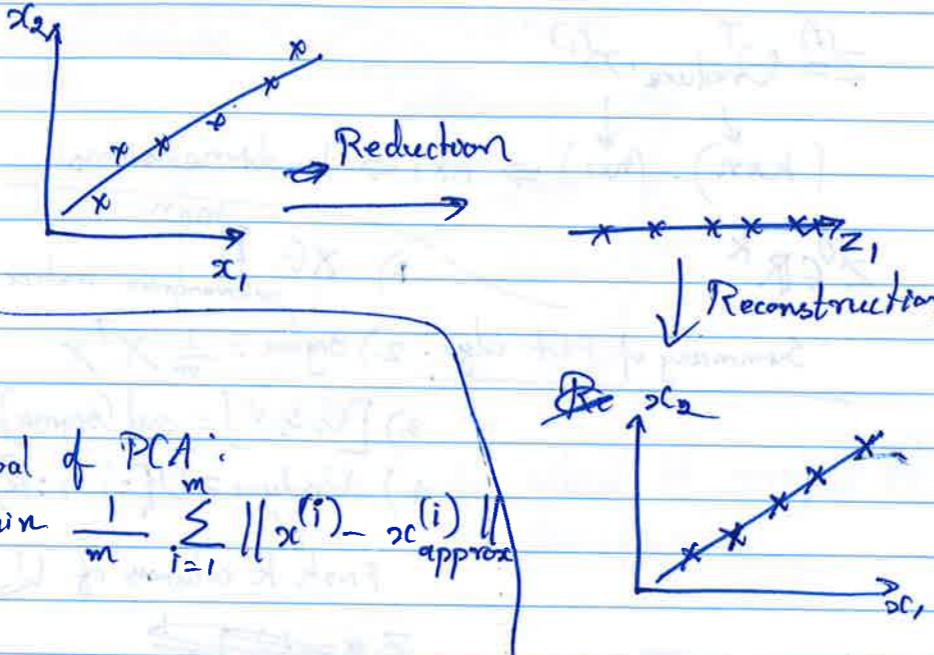
$$x_{\text{approx}} \approx x$$

$$x_{\text{approx}} = V_{\text{reduce}}^T \cdot z$$

$$X_{\text{approx}} = Z \cdot U^T$$

↓
m x n ↓
m x k k x n

Approx reconstruction of X



Goal of PCA:

$$\min \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|$$

Choosing the no. of principal components (k):

Choosing k :

PCA \rightarrow tries to minimize mean square projection error.

$$\left(\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2 \right)$$

$$\text{Total variation} = \frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$$

Covariance between $x_{\text{approx}}^{(i)}$ and origin

Choose the smallest k such that

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq 0.01 \quad (1\%)$$

→ 99% of variance is retained.

$$[U, S, V] = \text{svd}(\text{sigma})$$

$n \times n$ diagonal matrix.

$$S = \begin{bmatrix} s_{11} & 0 & 0 & 0 & 0 & \dots & 0 \\ 0 & s_{22} & 0 & \dots & 0 & \dots & 0 \\ \vdots & 0 & \ddots & \dots & \dots & \dots & 0 \\ 0 & 0 & \dots & \dots & \dots & \dots & s_{nn} \end{bmatrix}_{n \times n}$$

$$\frac{1 - \sum_{i=1}^k s_{ii}}{\sum_{i=1}^n s_{ii}} \leq 0.01$$

Summary of choosing k (no. of principal components):

$$[U, S, V] = \text{svd}(\text{sigma})$$

Pick smallest k for which

$$\frac{\sum_{i=1}^k s_{ii}}{\sum_{i=1}^n s_{ii}} \geq 0.99$$

(99% of variance is retained)

Applying PCA:

- Speedup supervised learning

$$(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$$

$$x^{(i)} \in \mathbb{R}^{10000}$$

\Rightarrow Slow learning (too many features)

Extract inputs: $x^{(1)}, \dots, x^{(m)}$

↓PCA

$$z^{(1)}, \dots, z^{(m)}$$

New training set:

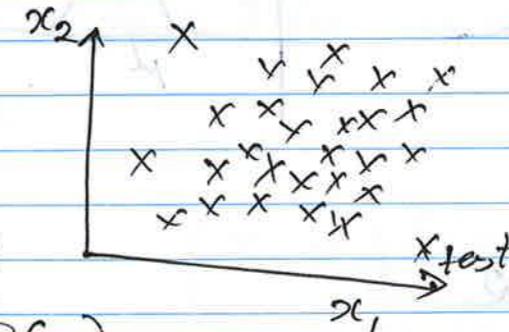
$$(z^{(1)}, y^{(1)}), \dots, (z^{(m)}, y^{(m)})$$

Note: Do NOT use PCA just to reduce overfitting

Week 9:

Anomaly detection:

Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$



Model

$$P(x)$$

$P(x_{test}) < \epsilon \downarrow$
Anomaly

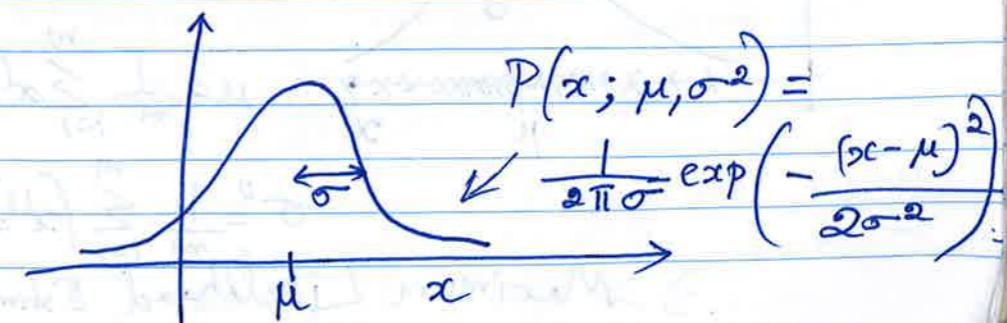
$P(x_{test}) \geq \epsilon \downarrow$
OK

Gaussian (Normal) distribution:

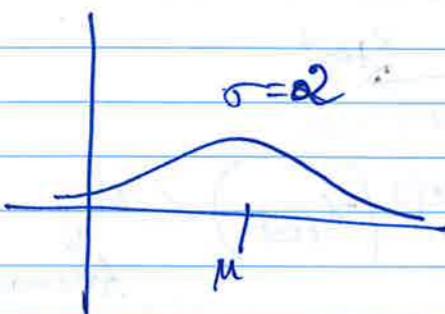
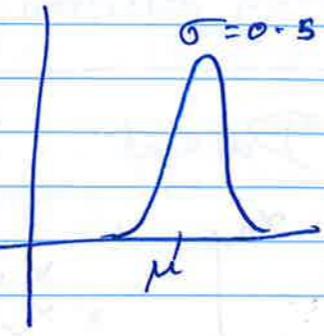
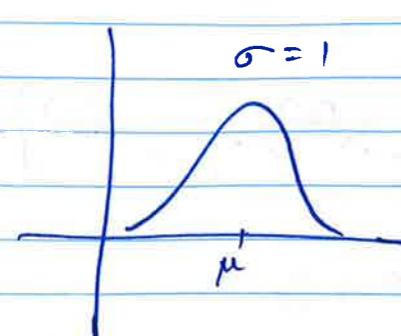
$$x \in \mathbb{R}, \quad x \sim N(\mu, \sigma^2)$$

"distributed as"

$\mu \rightarrow$ mean, $\sigma^2 \rightarrow$ variance

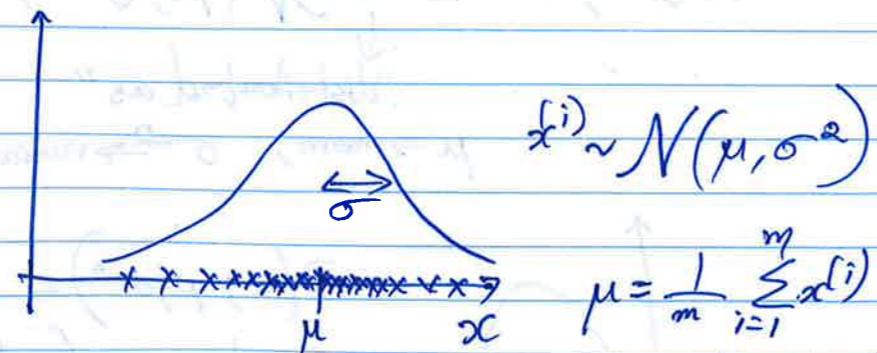


Examples:



Parameter estimation:

Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$,
 $x^{(i)} \in \mathbb{R}$



$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$
$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$

Maximum Likelihood Estimator

Density estimation:
(Anomaly detection)

→ Training set: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

→ For each example $x \in \mathbb{R}^n$

$$P(x) = P(x_1; \mu_1, \sigma_1^2) \cdot P(x_2; \mu_2, \sigma_2^2) \cdots P(x_n; \mu_n, \sigma_n^2)$$

$$x_1 \sim N(\mu_1, \sigma_1^2)$$

$$x_2 \sim N(\mu_2, \sigma_2^2)$$

$$\vdots$$
$$x_n \sim N(\mu_n, \sigma_n^2)$$

$$P(x) = \prod_{j=1}^n P(x_j; \mu_j, \sigma_j^2) \quad P(x) = \prod_{j=1}^n P(x_j; \mu_j, \sigma_j^2)$$

Anomaly detection algorithm:

1.) Choose features that might be indicative of anomalies

2.) Fit params $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j; \quad \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

3.) Given new example x , compute $P(x)$:

$$P(x) = \prod_{j=1}^n P(x_j; \mu_j, \sigma_j^2)$$
$$= \prod_{j=1}^n \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

Anomaly if $P(x) < \epsilon$.

Developing & Evaluating Anomaly Detection System:

Training set \rightarrow non-anomalous

Cross-validation \rightarrow few anomalous examples.
Test-set

\rightarrow Fit model $P(x)$ on training set

\rightarrow On cross-validation / test example x_i ,

$$y = \begin{cases} 1 & \text{if } P(x) < \epsilon \\ 0 & \text{if } P(x) \geq \epsilon \end{cases}$$

predict
(anomaly)
(normal)

Evaluation metrics:

- Precision / Recall
- F-score

Cross-validation set to choose param ϵ .

Anomaly detection vs Supervised learning:

Anomaly detection

- \rightarrow Very small number of positive examples,
- \rightarrow large number of negative examples.
- \rightarrow Many different types of anomalies

Supervised Learning

- Large no. of positive & negative examples.
- Sufficient positive examples for learning & prediction

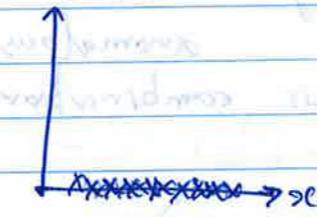
Choosing features for Anomaly detection:

Non-Gaussian features \rightarrow Transform them to see if we get Gaussian ($x \rightarrow \log(x)$ for example)

Error analysis:

Want large $P(x)$ for normal examples,
small $P(x)$ for anomalies.

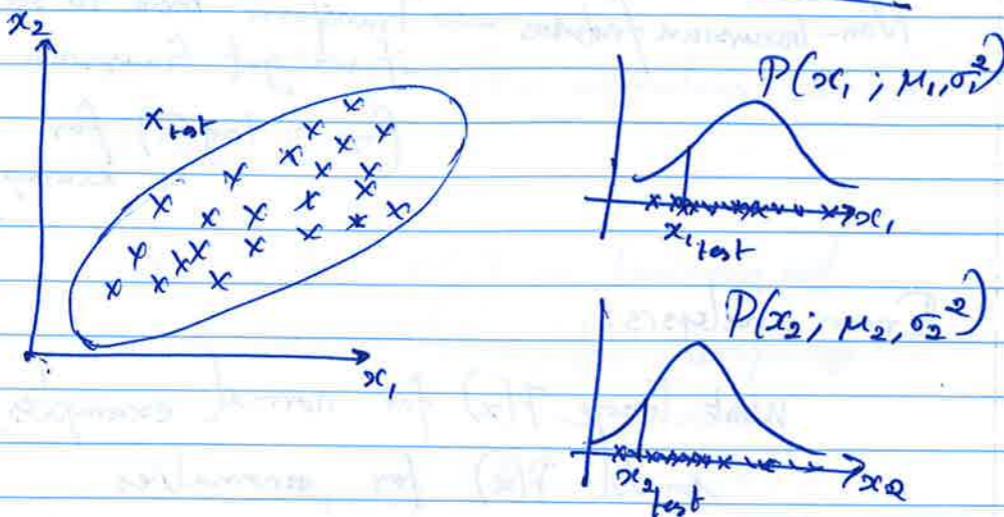
Problem: $P(x)$ is comparable for both



New features:

$$x_j = \frac{\text{CPU load}}{\text{Network traffic}}$$

Multivariate Gaussian Distribution:



Problem:

Individual features by themselves are not anomalous,
but their combination is

$$\downarrow \\ j=1^n P(x_j; \mu_j, \sigma_j^2)$$

fails to catch this

(each individual $P(x_j)$
will be high enough).

Solution: Multivariate Gaussian distribution

$x \in \mathbb{R}^n$. Don't model $P(x_1), P(x_2), \dots, P(x_n)$

Model $P(x)$ all in one go.

Params: $\mu \in \mathbb{R}^n$, $\Sigma \in \mathbb{R}^{n \times n}$
 \rightarrow Covariance Matrix

$$P(x; \mu, \Sigma) =$$

$$\frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

\rightarrow Determinant of Σ

$\Sigma \rightarrow$ Covariance matrix ($n \times n$)

1) Parameter fitting:

Training set: $\{x^{(1)}, \dots, x^{(m)}\}$

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

\downarrow
Covariance matrix

2) New example x :

$$P(x) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

$P(x) < \epsilon \Rightarrow$ anomaly

Original model (Univariate Gaussian)

$$P(x) = \prod_{j=1}^n P(x_j; \mu_j, \sigma_j^2)$$



This corresponds to multivariate

Gaussian $P(x; \mu, \Sigma)$

where

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & 0 & \dots & 0 \\ 0 & 0 & \ddots & \ddots & 0 \\ \vdots & & & \ddots & 0 \\ 0 & 0 & \dots & 0 & \sigma_n^2 \end{bmatrix}$$

Covariance matrix

Meaning, features
are linearly independent

~~or~~
Univariate Gaussian

- Manually create features to capture relationship between features

- Computationally cheaper

- Okay if m (training set size) is small.

Multivariate Gaussian

Automatically captures correlation b/w features.

Computationally expensive
(computing Σ^{-1})

- $m > n$, otherwise
 Σ is non-invertible

Recommender Systems:

→ Content-based recommendation:

→ Create features for each item to predict

→ Create feature vector

→ For each user, learn params $\alpha^{(j)} \in \mathbb{R}^3$

→ Prediction for user j :

$$(\alpha^{(j)})^T \cdot x$$

$r(i, j) = 1$ if user j has rated movie i

$y^{(i,j)}$ = rating by user j on movie i

$\alpha^{(j)}$ = parameter vector for user j

$x^{(i)}$ = feature vector for movie i

For user j , movie i , predicted rating:

$$(\alpha^{(j)})^T \cdot x^{(i)}$$

$m^{(j)}$ = no. of movies rated by user j

To learn $\alpha^{(j)}$:

$$\min_{\alpha^{(j)}} \frac{1}{2m^{(j)}} \sum_{i: r(i,j)=1} \left((\alpha^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2m^{(j)}} \sum_{k=1}^{m^{(j)}} \alpha_k^{(j)^2}$$

(Linear regression)

To learn $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(nu)}$:

$$\min_{\theta^{(1)}, \dots, \theta^{(nu)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i: r(i,j)=1} \left((\theta^{(i)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

$J(\theta^{(1)}, \dots, \theta^{(nu)})$

Gradient descent:

$$\theta_k^{(j)} = \theta_k^{(j)} - \alpha \sum_{i: r(i,j)=1} \left((\theta^{(i)})^T x^{(i)} - y^{(i,j)} \right) x_k^{(i)} \quad (\text{for } k=0)$$

$$\theta_k^{(j)} = \theta_k^{(j)} - \alpha \left(\sum_{i: r(i,j)=1} \left((\theta^{(i)})^T x^{(i)} - y^{(i,j)} \right) x_k^{(i)} + \lambda \theta_k^{(j)} \right) \quad (\text{for } k \neq 0)$$

Collaborative Filtering: → Learn features for item x

Given $\theta^{(1)}, \dots, \theta^{(nu)}$, to learn $x^{(i)}$:

$$\min_{x^{(i)}} \frac{1}{2} \sum_{j: r(i,j)=1} \left((\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{k=1}^n (x_k^{(i)})^2$$

Given $\theta^{(1)}, \dots, \theta^{(nu)}$, to learn $x^{(1)}, \dots, x^{(nm)}$:

$$\min_{x^{(1)}, \dots, x^{(nm)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j: r(i,j)=1} \left((\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2$$

Collaborative Filtering Algorithm:

1. Given $x^{(1)}, \dots, x^{(nm)}$, estimate $\theta^{(1)}, \dots, \theta^{(nu)}$
2. Given $\theta^{(1)}, \dots, \theta^{(nu)}$, estimate $x^{(1)}, \dots, x^{(nm)}$
3. Repeat (1) & (2) till convergence.

(Or)
Minimizing $x^{(1)}, \dots, x^{(nm)}$ and $\theta^{(1)}, \dots, \theta^{(nu)}$ simultaneously

$$J(x^{(1)}, \dots, x^{(nu)}, \theta^{(1)}, \dots, \theta^{(nu)}) = \frac{1}{2} \sum_{(i,j): r(i,j)=1} \left((\theta^{(j)})^T x^{(i)} - y^{(i,j)} \right)^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

$$\min_{x^{(1)}, \dots, x^{(nu)}, \theta^{(1)}, \dots, \theta^{(nu)}} J(x^{(1)}, \dots, x^{(nu)}, \theta^{(1)}, \dots, \theta^{(nu)})$$

Collaborative filtering algo:

- 1) Initialize $x^{(1)}, \dots, x^{(nm)}, \theta^{(1)}, \dots, \theta^{(nu)}$ to small random values.
- 2) Minimize $J(x^{(1)}, \dots, x^{(nm)}, \theta^{(1)}, \dots, \theta^{(nu)})$ using gradient descent (or other optimization).
- 3) For a user with params θ and a movie with learned features x , predict rating as $\theta^T x$.

Vectorization: Low Rank Matrix Factorization

$$Y = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$$

such that $y^{(i,j)} = \text{rating of movie } i \text{ by user } j$

Predicted ratings:

$$\begin{array}{c} \text{User 1} \quad \text{User 2} \quad \dots \quad \text{User } n_u \\ \text{Movie 1} \quad (\theta^{(1)})^T (x^{(1)}) \quad (\theta^{(2)})^T (x^{(1)}) \quad \dots \quad (\theta^{(nu)})^T (x^{(1)}) \\ \text{Movie 2} \quad (\theta^{(1)})^T (x^{(2)}) \quad ((\theta^{(2)})^T (x^{(2)})) \quad \dots \quad (\theta^{(nu)})^T (x^{(2)}) \\ \vdots \\ \text{Movie } nm \quad (\theta^{(1)})^T (x^{(nm)}) \quad (\theta^{(2)})^T (x^{(nm)}) \quad \dots \quad ((\theta^{(nu)})^T (x^{(nm)})) \end{array}$$

$$X = \begin{bmatrix} - (x^{(1)})^T - \\ - (x^{(2)})^T - \\ \vdots \\ - (x^{(nm)})^T - \end{bmatrix}; \theta = \begin{bmatrix} - (\theta^{(1)})^T \\ - (\theta^{(2)})^T \\ \vdots \\ - (\theta^{(nu)})^T \end{bmatrix}$$

$$\text{Predicted ratings} = X \theta^T$$

↓
Low-rank matrix factorization.

Mean Normalization:

→ If a user hasn't rated any movie,

$$y^{(i,j)} = ? \text{ for all movies } i.$$

↳ $(\theta^{(i)})^T x^{(i)} - y^{(i,j)} \rangle^2$ term of cost fn., is zero

↑
Not useful.

→ mean normalize the ratings for each movie.

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_{nm} \end{bmatrix}, \mu_i \rightarrow \text{mean rating of movie } i$$

For user j on movie i , predict
 $(\theta^{(j)})^T (x^{(i)}) + \mu_i$

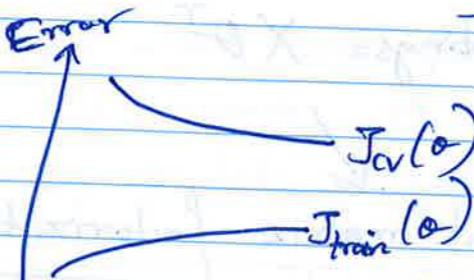
Week 10:

Learning with large datasets:

→ First check if ~~large~~
more training examples

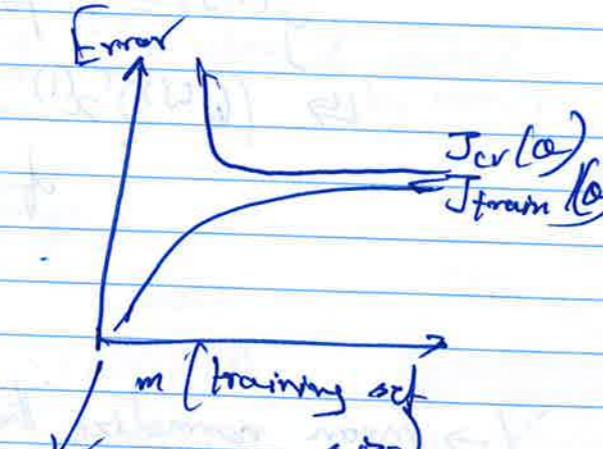
will help with learning

→ Plot the learning curve & verify



Overfits when
(because) m is small

More data will
help



Underfits even with small m
(High bias)

Increasing m won't help

Stochastic Gradient Descent (SGD):

Linear regression gradient descent:

$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j = \theta^T x$$

$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m \left[(h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \right]$$

for $j = 0, \dots, n$

For large m, computing $\frac{\partial J(\theta)}{\partial \theta_j}$ is expensive

→ Batch gradient descent

$$\text{SGD: cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$J_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

1. Randomly shuffle dataset.

2. Repeat {

for $i = 1, \dots, m$ {

$$\theta_j := \theta_j - \alpha \left(h_{\theta}(x^{(i)}) - y^{(i)} \right) \cdot x_j^{(i)}$$

$$\left(\text{for } j = 0, \dots, n \right)$$

$$\frac{\partial}{\partial \theta_j} \text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

B3

Mini-batch Gradient Descent:

Batch Gradient Descent \rightarrow Use all m examples for each iteration

Stochastic Gradient Descent \rightarrow Use 1 example in each iteration

Minibatch Gradient Descent \rightarrow Use b examples in each iteration
($b = \text{minibatch size}$)

Say $b = 10$, $m = 1000$

Repeat {

$$\text{for } i = 1, 11, 21, 31, \dots, 991 \quad \{$$

$$\theta_j := \theta_j - \alpha \cdot \frac{1}{10} \sum_{k=i}^{i+9} \left(h_\theta(x^{(k)}) - y^{(k)} \right) x_j^{(k)}$$

(for every $j = 0, \dots, n$)

}

\rightarrow Better vectorization over SGD.

\rightarrow But, need to fit another param 'b'.

(minibatch size)

SGD convergence check:

Batch gradient descent convergence check:

\rightarrow Plot $J(\theta)$ as a function

\rightarrow Plot $J_{\text{train}}(\theta)$ vs no. of iterations,

\rightarrow will always decrease (non-increasing) with each iteration.

Note $\rightarrow J(\theta)$ can increase after each iteration.

\hookrightarrow Checking for convergence:

$$\rightarrow \text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

$$= \frac{1}{2} \sum_{j=0}^n (h_\theta(x^{(i)}) - y^{(i)})^2$$

\rightarrow During learning, compute $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ before updating θ using $(x^{(i)}, y^{(i)})$.

\rightarrow Every 1000 iterations, plot $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ averaged over the last 1000 examples.

Learning rate in SGD:

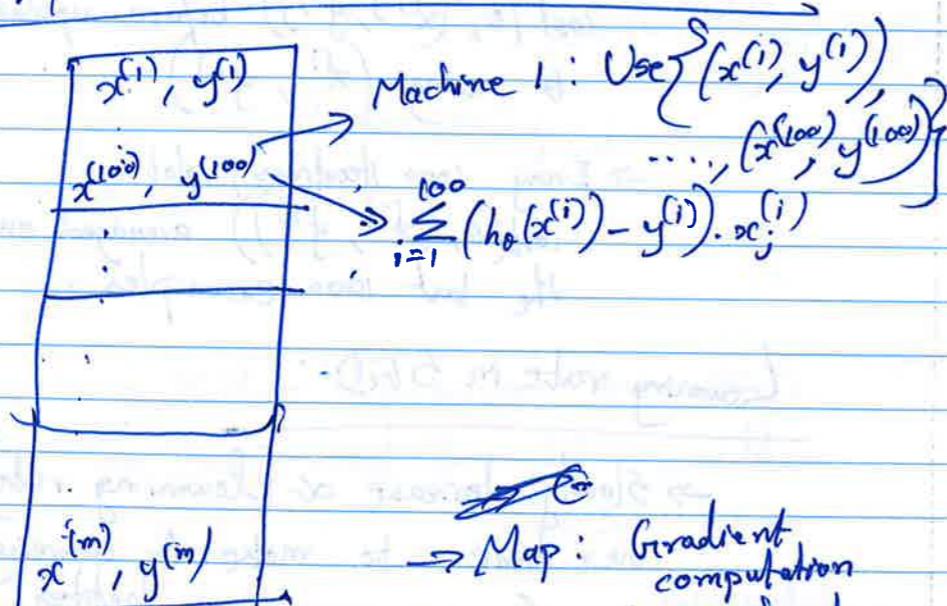
\rightarrow Slowly decrease α (learning rate) over time to make θ converge better.

$$\left(\text{e.g., } \alpha = \frac{\text{const}_1}{(\text{Iteration Number} + \text{const}_2)} \right)$$

Online Learning:

- One example at a time
- Use it to update params, and then discard the example (just one iteration or update, example never used again)
- ~~Adaptive w~~
- No fixed training set.
- Can adapt to changes in ~~the~~ model over time.

Map reduce & Data parallelism:



→ Map: Gradient computation for subset of examples.

→ Reduce: Average all gradients, update params.

Week 11

Photo OCR:

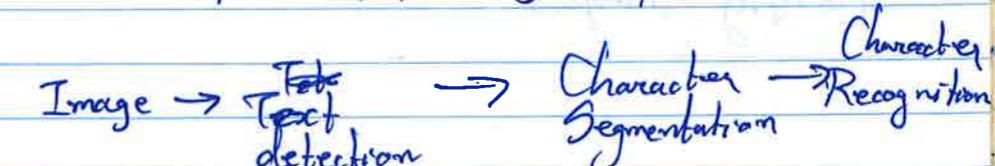
Pipeline:

Step 1: Text detection

Step 2: Character segmentation

Step 3: Character classification

Step 4: (Optional) Spell correction



Sliding Windows:

Text detection train.

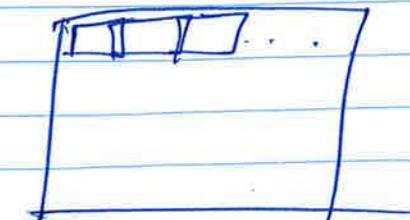
→ ~~Text~~ Image patches with and without text



→ Sliding window detection for text

- Sliding window over the entire image to get the patch

- Stride



→ 1D sliding window for
character segmentation ~~segmentation~~

→ Classify characters

Artificial Data Synthesis:

Ceiling Analysis: